

Bayesian Methods for Genetic Association Analysis with Heterogeneous Subgroups: from Meta-Analyses to Gene-Environment Interactions

Xiaoquan Wen^{*1} and Matthew Stephens^{†2,3}

¹Department of Biostatistics, School of Public Health, University of Michigan

²Department of Statistics, University of Chicago

³Department of Human Genetics, University of Chicago

November 9, 2011

Abstract

In genetic association analyses, it is often desired to analyze data from multiple potentially-heterogeneous subgroups. The amount of expected heterogeneity can vary from modest (as might typically be expected in a meta-analysis of multiple studies of the same phenotype, for example), to large (e.g. a strong gene-environment interaction, where the environmental exposure defines discrete subgroups). Here, we consider a flexible set of Bayesian models and priors that can capture these different levels of heterogeneity. We provide accurate numerical approaches to compute approximate Bayes Factors for these different models, and also some simple analytic forms which have natural interpretations and, in some cases, close connections with standard frequentist test statistics. These approximations also have the convenient feature that they require only summary-level data from each subgroup (in the simplest case, a point estimate for the genetic effect, and its standard error, from each subgroup). We illustrate the flexibility of these approaches on three examples: an analysis of a potential gene-environment interaction for a recombination phenotype, a large scale meta-analysis of genome-wide association data from the Global Lipids consortium, and a cross-population analysis for expression quantitative trait loci (eQTLs).

^{*}xwen@umich.edu

[†]mstephens@uchicago.edu

1 Introduction

In this paper, we develop Bayesian methods for analyzing genetic association data, allowing for potential heterogeneity among (pre-specified) subgroups. We are motivated by two distinct settings where heterogeneity may arise. The first setting is meta-analysis of multiple association studies of the same phenotype. These studies are usually carried out by different investigators, at different centers, and so heterogeneity (or apparent heterogeneity) of genetic effects might be expected (e.g. due to differences in the way phenotypes are measured, or due to systematic differences between individuals enrolled in each study). Such meta-analyses have become an increasingly popular and important statistical tool for detecting modest genetic associations that are too small to be detected in smaller individual studies (Teslovich et al. (2010), Zeggini et al. (2008)), and frequentist methods focused on allowing for heterogeneity in this setting were recently proposed by Lebrech et al. (2010) and Han and Eskin (2011). The second setting is where genuine biological interactions may cause some genetic variants to exhibit different effects on individuals in different subgroups; for example, genetic effects can differ in males and females even at autosomal loci (Kong et al. (2008), Ober et al. (2008)). And in gene expression analyses that aim to detect genetic variants associated with gene expression levels, data are often available on individuals from different continental groups (Stranger et al. (2007), Veyrieras et al. (2008)), or on different tissue types (Dimas et al. (2009)), where heterogeneity of effects may be expected.

These two settings differ in the extent of the heterogeneity expected: for example, interactions could cause genetic variants to have effects in different directions in different subgroups, whereas this might be considered unlikely in the meta-analysis setting. They also differ in the extent to which heterogeneity may be of direct interest (e.g. in interactions) or largely a “nuisance” (e.g. in meta-analysis). However, the two settings also share an important element in common: the vast majority of genetic variants are unassociated with any given phenotype of interest, within *all* subgroups. Consequently, it is of considerable interest to identify genetic variants that show association in *any* subgroup, or in other words to reject the “global” null hypothesis of *no association within any subgroup*. This focus on rejecting the global null hypothesis distinguishes genetic association analyses from other settings, and calls for analysis approaches tailored to this goal; see Lebrech et al. (2010) for relevant discussion.

The methods we consider here are suited to the analysis of genetic association studies, and are sufficiently flexible to handle a range of settings in which heterogeneity may be an issue, from meta-analyses to gene-environment interactions. In brief, we introduce families of alternative models that allow for a range of different effect sizes and levels of heterogeneity, and then address questions of interest by comparing the support in the data for these different models vs the global null. Within this framework the goal of testing the global null is accomplished by assessing the overall support for any of the alternative hypotheses, whereas the goal of examining heterogeneity among groups is achieved by comparing the relative support for different alternative models.

Although there has been previous work on Bayesian methods for meta-analysis (e.g. Sutton and Abrams (2001), Stangl and Berry (2000), DuMouchel and Harris (1983), Whitehead and Whitehead (1991), Li and Begg (1994), Eddy et al. (1990), Givens et al. (1997), Verzilli et al. (2008), De Iorio et al. (2011), Burgess et al. (2010), Mila and Ngugi (2011)), our focus here is somewhat different to much of this. First, due to the nature of genetic association studies mentioned above, our focus is particularly on testing and model comparison, via Bayes Factors, rather than on estimation. Second, due to the fact that in genetic studies one often wishes to perform millions of analyses, we focus on obtaining fast numerical approximations to Bayes Factors. Finally, we provide novel results that connect the Bayesian methods with standard frequentist test statistics used in many meta-analyses. Specifically we show how for certain prior assumptions the Bayesian analyses depend on the phenotype data, asymptotically, only through these test statistics, and that the common practice of ranking SNP associations by these standard test statistics is equivalent to making specific (and not necessarily realistic) prior assumptions.

2 Models and Methods

We start with specifying models for quantitative traits, derive efficient methods for computing (approximate) Bayes Factors comparing these models, and discuss their statistical properties and connections with frequentist test statistics. Later, we generalize these results to binary outcomes in case-control studies.

2.1 Notation and Assumptions

Assume (quantitative) phenotype data and genotype data are available on S pre-defined subgroups, and focus on assessing association between the phenotype and each genetic variant, one at a time. For convenience we assume each genetic variant is a Single Nucleotide Polymorphism (SNP), although our methods could easily accommodate other types of variant. We assume that the data within subgroup s come from n_s randomly-sampled unrelated individuals. Let the n_s -vectors \mathbf{y}_s and \mathbf{g}_s denote, respectively, the corresponding phenotype data and the genotype data at a single “target” SNP. We also let $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_S)$ and $\mathbf{G} = (\mathbf{g}_1, \dots, \mathbf{g}_S)$ denote the complete set of phenotype and genotype data respectively.

2.2 Hierarchical Models for Quantitative Traits

In this section, we introduce a set of models for describing heterogeneous genetic effects for a target SNP across subgroups.

Within each subgroup, we model the association between phenotype and genotype using a standard linear model. Without loss of generality, in subgroup s , we assume

$$\mathbf{y}_s = \mu_s \mathbf{1} + \beta_s \mathbf{g}_s + \mathbf{e}_s, \quad \mathbf{e}_s \sim N(0, \sigma_s^2 I). \quad (2.1)$$

Here, we also assume residual errors are independent across subgroups.

The “global” null hypothesis of interest is that there is no genotype-phenotype association within any subgroup; that is, $\beta_s = 0$ for all s .

Under the alternative hypothesis we begin by assuming that the genetic effects among subgroups are *exchangeable*, and more specifically that they are normally distributed about some unknown common mean. We consider two different definitions of genetic effects: the “standardized effects” $b_s := \beta_s / \sigma_s$, and the unstandardized effects, β_s , leading to the following models:

1. *Exchangeable Standardized Effects (ES model)*. Under this model we assume that the standardized

effects b_s are normally distributed among subgroups:

$$b_s|\sigma_s \sim N(\bar{b}, \phi^2) \text{ or equivalently, } \beta_s|\sigma_s \sim N(\sigma_s \bar{b}, \sigma_s^2 \phi^2), \quad (2.2)$$

so the hyper-parameters \bar{b} and ϕ characterize, respectively, the mean and variance of effects among subgroups. We also assume a normal prior distribution for \bar{b} ,

$$\bar{b} \sim N(0, \omega^2). \quad (2.3)$$

Finally, for the parameters σ_s and μ_s , which are common to both the null and alternative hypotheses, we use the convenient conjugate priors

$$\mu_s|\sigma_s \sim N(0, \sigma_s^2 v_s^2); \quad \sigma_s^{-2} \sim \Gamma(m_s/2, l_s/2). \quad (2.4)$$

When performing inference we consider the posterior distributions that arise in the limits $v_s^2 \rightarrow \infty$ and $l_s, m_s \rightarrow 0$ (which correspond to standard improper priors for a normal mean and variance, and result in proper posteriors).

2. *Exchangeable Effects (EE model)*. Under this model we assume that the unstandardized effects β_s are normally distributed:

$$\beta_s \sim N(\bar{\beta}, \psi^2), \quad (2.5)$$

where $\bar{\beta}$ and ψ play similar roles to \bar{b} and ϕ in the ES model. We also assume a normal prior for $\bar{\beta}$,

$$\bar{\beta} \sim N(0, w^2), \quad (2.6)$$

and priors for (μ_s, σ_s) :

$$\mu_s \sim N(0, u_s^2); \quad \sigma_s^{-2} \sim \Gamma(m_s/2, l_s/2). \quad (2.7)$$

For each subgroup this prior specification is effectively the semi-conjugate prior commonly used in Bayesian linear regression. As above, we consider the limits $u_s^2 \rightarrow \infty$ and $l_s, m_s \rightarrow 0$.

In both the ES and EE models the alternative hypothesis involves two key hyper-parameters, one (ω in the ES model and w in the EE model) that controls the prior expected size of the average effect across subgroups, and another (ϕ in the ES model and ψ in the EE model) that controls the prior expected degree of heterogeneity among subgroups. A complimentary view is that $\omega^2 + \phi^2$ (respectively, $w^2 + \psi^2$) controls the expected (marginal) effect size in each study and ϕ/ω (respectively, ψ/w) controls the degree of heterogeneity.

Of the two models, the ES model has the advantage that it results in analyses (e.g. Bayes Factors) that are invariant to the phenotype measurement scale used within each subgroup. This not only makes it more robust to users accidentally specifying phenotype measurements in different subgroups on different scales (possibly a non-trivial issue in complex analyses involving collaboration among many research groups), but also means that it can be applied when measurement scales may be difficult to harmonize across subgroups, for example due to the use of different measurement technologies. For these reasons we prefer the ES model for general use. However, in some cases the EE model may be easier to apply. For example, if one has access only to published point estimates and standard errors for the effect size β_s in each study, then this suffices to approximate the Bayes Factor under the EE model, but not under the ES model. Note that the ES and EE models will produce similar results to each other if the residual error variances are similar in all subgroups.

It is also possible to control for additional (possibly study-specific) covariates by including additional predictor variables into (2.1). If independent flat priors are used for the coefficients of these controlled covariates within each study then our main results below still hold, effectively unchanged. This treatment is analogous to the frequentist mixed effects model, where controlled covariates are typically assumed to have study-specific effects.

2.2.1 A Curved Exponential Family Normal Prior

In the priors described above the variance of the effect sizes (ω^2) is independent of the average effect (\bar{b}). In some settings this independence assumption may seem unattractive. For example, in meta-analysis, we may expect genuine genetic association to possess the property that effect sizes across studies typically have the same sign (Owen (2009)), regardless whether \bar{b} is small or large. But the independence assumption implies that the probability that the effects have the same sign is much larger when \bar{b} is large than when it is small. To address this we can modify the priors above to allow the variance to depend on the mean. For example, under the ES model, we can replace (2.2) with

$$b_s \sim N(\bar{b}, k^2 \bar{b}^2), \quad (2.8)$$

so that

$$\Pr(b_s \text{ has a different sign from } \bar{b}) = \Phi\left(-\frac{1}{|k|}\right), \quad (2.9)$$

where Φ is the cumulative probability function of standard normal distribution. Note that (2.9) does not depend on \bar{b} . For example, when $k = 1/2$, sampling from this prior distribution, the probability of obtaining a value of b_s having an opposite sign to \bar{b} is approximately 2.3%. As the value of k decreases, the restriction becomes more stringent. When $k = 0$, the prior indicates all b_s are exactly same as \bar{b} , i.e. there is no heterogeneity of effects across subgroups.

Similarly, for the EE model, we can replace (2.5) with

$$\beta_s \sim N(\bar{\beta}, k^2 \bar{\beta}^2). \quad (2.10)$$

We refer to these alternative priors as “Curved Exponential Family Normal” (CEFN) priors, because they involve a functional relationship between the mean and variance.

2.3 Bayes Factors for Testing the Global Null Hypothesis

We now consider computing Bayes Factors for testing the “global” null hypothesis that the phenotype is not associated with the target SNP in any of the subgroups, versus the alternative hypotheses outlined above. We focus primarily on the simplest ES model, and give details for the EE model, and modifications for the CEFN models, in appendices.

Recall that the ES model is indexed by two parameters, ϕ and ω . Within this model, the global null hypothesis, which is most naturally written as $\beta_s \equiv 0$ for all s , can also be written as

$$H_0 : \phi = \omega = 0. \quad (2.11)$$

To compare the support in the data for this null hypothesis with the support for a particular alternative ES model specified by parameters (ϕ, ω) , we use the Bayes Factor:

$$\text{BF}^{\text{ES}}(\phi, \omega) = \frac{P(\mathbf{Y}|\mathbf{G}, \phi, \omega)}{P(\mathbf{Y}|\mathbf{G}, H_0)}. \quad (2.12)$$

This Bayes Factor also depends on the hyper-parameters v_s, l_s and m_s ; however, because these hyper-parameters are common to both the null and alternative hypotheses, the value of the Bayes Factor is not especially sensitive to the values chosen. As noted above we take the limits

$$v_s^2 \rightarrow \infty, \quad l_s \rightarrow 0, \quad m_s \rightarrow 0, \quad \forall s. \quad (2.13)$$

Each value of (ω, ϕ) corresponds to a particular alternative model, with ω controlling the typical average effect size, and ϕ controlling the degree of heterogeneity among subgroups (or in a re-parametrization, $\omega^2 + \phi^2$ controls the expected marginal effect size in each subgroup and ϕ/ω controls the degree of heterogeneity). There may be reasonable uncertainty about appropriate values for ϕ and ω due to the unknown mechanisms that cause heterogeneity. One simple way to allow for this is to specify a (discrete) prior distribution on a set of plausible values $\{(\phi^{(i)}, \omega^{(i)}) : i = 1, \dots, M\}$. We give a specific choice of such prior in the applications below. If π_i denotes the prior weight on $(\phi^{(i)}, \omega^{(i)})$ then the

resulting Bayes Factor against H_0 is the weighted average of the individual BFs:

$$\text{BF}_{\text{av}}^{\text{ES}} := \sum_{i=1}^M \pi_i \text{BF}^{\text{ES}}(\phi^{(i)}, \omega^{(i)}). \quad (2.14)$$

This average could also, of course, be extended to include averages over other models (e.g. one could average over some models that use the CEFN prior, and others that do not). The fact that Bayes Factors under different models for heterogeneity can be both averaged in this way (to assess evidence against the global null, allowing for heterogeneity), and compared with one another (to assess the evidence for different levels or types of heterogeneity), is one advantage of the Bayesian framework compared with typical frequentist treatments.

2.3.1 Calculating the Bayes Factors

Calculating $\text{BF}^{\text{ES}}(\phi, \omega)$ boils down to evaluating a complicated multi-dimensional integral. In appendix A we present two different approximations, both based on applying Laplace’s method and both having error terms that decay inversely with the average sample size across subgroups. The first of these, which effectively follows methods from Butler and Wood (2002) for computing confluent hyper-geometric functions, is very accurate, even for small sample sizes. Indeed, for the special case of a single subgroup ($S = 1$), the approximation becomes *exact*, and for small numbers of subgroups we have checked numerically (appendix D) that it provides almost identical results to an alternative approach based on adaptive quadrature (which is practical only for small S). However, it requires a numerical optimization step and has a somewhat complex form, which although not a practical barrier to its use does hinder intuitive interpretation. In what follows we use $\widehat{\text{BF}}^{\text{ES}}$ to denote this approximation.

The second approximation is less accurate for small samples sizes, but converges asymptotically (with average sample size) to the correct answer. For the special case of $S = 1$ it yields an analogue of the approximate Bayes Factors from Wakefield (2009) and Johnson (2008), and in what follows we use ABF^{ES} to denote this approximation under the ES model. The nice feature of ABF^{ES} is that it has an intuitive analytic form with close connections to standard frequentist test statistics for meta-analysis. Proposition 1 below gives this analytic form in detail.

Before stating Proposition 1, we introduce some notation.

- *Association Testing in a Single Subgroup*

First, consider analyzing a single subgroup, s . Let $\hat{\beta}_s$ and $\hat{\sigma}_s$ denote the least square estimates of β_s and σ_s from the linear regression model (2.1) using only data from s . Then an estimate for the standardized effect b_s can be obtained from

$$\hat{b}_s = \hat{\beta}_s / \hat{\sigma}_s. \quad (2.15)$$

Under the assumption of no association, the standard error of \hat{b}_s , $\delta_s := \text{se}(\hat{b}_s)$, is

$$\delta_s^2 = \frac{1}{\mathbf{g}'_s \mathbf{g}_s - n_s \bar{g}_s^2}. \quad (2.16)$$

A statistic T_s for testing $b_s = 0$ can be constructed as

$$T_s^2 = \frac{\hat{b}_s^2}{\text{se}(\hat{b}_s)^2} = \frac{\hat{\beta}_s^2}{\hat{\sigma}_s^2 \delta_s^2}. \quad (2.17)$$

Note that T_s is also equal to $\hat{\beta}_s / \text{se}(\hat{\beta}_s)$, which is the usual t-statistic for testing $\beta_s = 0$ in a frequentist framework.

Both Wakefield (2009) and Johnson (2008) derive an approximate Bayes Factor for testing $b_s \sim N(0, \phi^2)$ vs. $b_s = 0$, which has the form

$$\text{ABF}_{\text{single}}^{\text{ES}}(T_s, \delta_s; \phi) := \sqrt{\frac{\delta_s^2}{\delta_s^2 + \phi^2}} \exp\left(\frac{T_s^2}{2} \frac{\phi^2}{\delta_s^2 + \phi^2}\right). \quad (2.18)$$

As noted by Wakefield, if ϕ is chosen differently for each SNP, and proportional to the value of δ_s^2 for that SNP, then $\text{ABF}_{\text{single}}^{\text{ES}}$ will rank the SNPs in the same way as the usual test statistic T_s . This result connects the standard frequentist analysis to a particular (approximate) Bayesian analysis in the case of a single subgroup. Proposition 1 below effectively extends this to multiple subgroups, allowing for heterogeneity among subgroups.

- *Testing Average Effect in a Random-effect Meta-analysis Model*

Now consider the standard frequentist test of $\bar{b} = 0$ in a random-effect meta-analysis of all subgroups, with $b_s \sim N(\bar{b}, \phi^2)$. If ϕ is considered known, the standard estimate for \bar{b} and its standard error $\zeta := \text{se}(\hat{\bar{b}})$ are

$$\hat{\bar{b}} = \frac{\sum_s (\delta_s^2 + \phi^2)^{-1} \hat{b}_s}{\sum_s (\delta_s^2 + \phi^2)^{-1}}, \quad (2.19)$$

and

$$\zeta^2 = \frac{1}{\sum_s (\delta_s^2 + \phi^2)^{-1}}. \quad (2.20)$$

The usual frequentist statistic $\mathcal{T}_{\text{ES}}^2$ for testing $\bar{b} = 0$ is

$$\mathcal{T}_{\text{ES}}^2 = \frac{\hat{\bar{b}}^2}{\text{se}(\hat{\bar{b}})^2}. \quad (2.21)$$

Applying Johnson's idea (Johnson (2005, 2008)) we can “translate” this test statistic into an approximate Bayes Factor for testing $\bar{b} \sim N(0, \omega^2)$ vs. $\bar{b} = 0$, which yields

$$\text{ABF}_{\text{single}}^{\text{ES}}(\mathcal{T}_{\text{ES}}^2, \zeta; \omega) := \sqrt{\frac{\zeta^2}{\zeta^2 + \omega^2}} \exp\left(\frac{\mathcal{T}_{\text{ES}}^2}{2} \frac{\omega^2}{\zeta^2 + \omega^2}\right). \quad (2.22)$$

With these ideas and notations in hand we can now describe the analytic form of the overall approximate Bayes Factor $\text{ABF}^{\text{ES}}(\phi, \omega)$, as a simple product of the ABFs (2.18) and (2.22).

PROPOSITION 1. *Under the ES model, applying a version of Laplace's method to approximate the Bayes Factor $\text{BF}^{\text{ES}}(\phi, \omega)$ yields the approximation*

$$\text{BF}^{\text{ES}}(\phi, \omega) \approx \text{ABF}^{\text{ES}}(\phi, \omega) := \text{ABF}_{\text{single}}^{\text{ES}}(\mathcal{T}_{\text{ES}}^2, \zeta; \omega) \cdot \prod_s \text{ABF}_{\text{single}}^{\text{ES}}(T_s^2, \delta_s; \phi), \quad (2.23)$$

and $\text{ABF}^{\text{ES}}(\phi, \omega)$ converges (almost surely) to $\text{BF}^{\text{ES}}(\phi, \omega)$ as $n_s \rightarrow \infty$ for all subgroups s .

Proof. See appendix A.1. □

NOTE 1. *If the study-specific residual error variances, σ_s are considered known (rather than being*

assigned a prior distribution), and used in place of $\hat{\sigma}_s$ in the calculations of ABF^{ES} , then the approximation is exact, and $\text{ABF}^{\text{ES}}(\phi, \omega) = \text{BF}^{\text{ES}}(\phi, \omega)$. This fact, together with the fact that the estimators $\hat{\sigma}_s$ are consistent for σ_s , provides an intuitive explanation for why the proposition holds.

NOTE 2. The numerical accuracy of ABF^{ES} as an approximation to BF^{ES} depends on sample sizes, and for small sample sizes it may not be sufficiently accurate for routine application. However, a simple modification, described in appendix C, yields much greater accuracy.

Proposition 1 breaks down the overall evidence for association into parts that are due to the evidence in each individual subgroup (the second term) and a part that reflects the consistency of the effects across subgroups (the first term). In particular, if all subgroups show effects in the same direction, then the first term will tend to be large ($\gg 1$) and provide a “boost” in the evidence for association compared with the situations when the effects across subgroups are in different directions.

A similar result holds for the EE model and is given in appendix A.2. The detailed computation of Bayes Factors involving CEFN priors is shown in appendix A.3. In appendix D we compare the numerical accuracies of the various Bayes Factor approximations described above.

2.4 Properties of Bayes Factors

In this section, we discuss some interesting and important properties of the Bayes Factors described above.

2.4.1 Bayes Factors Depend Only on Summary Statistics

Both the true Bayes Factors ($\text{BF}^{\text{ES}}, \text{BF}^{\text{EE}}$) and the approximations, $(\widehat{\text{BF}}^{\text{ES}}, \widehat{\text{BF}}^{\text{EE}}, \text{ABF}^{\text{ES}}, \text{ABF}^{\text{EE}})$ depend on the observed data in each subgroup only through a set of summary statistics, i.e., a 6-tuple $(n_s, \mathbf{1}'\mathbf{y}_s, \mathbf{1}'\mathbf{g}_s, \mathbf{y}_s'\mathbf{y}_s, \mathbf{g}_s'\mathbf{g}_s, \mathbf{y}_s'\mathbf{g}_s)$. Therefore, to use the proposed Bayesian methods it is not necessary to have access to the full data set from each individual subgroup. Further, for computing the ABFs the summary statistics needed from each subgroup are reduced to only $(\hat{b}_s, \text{se}(\hat{b}_s))$ for the ES model and $(\hat{\beta}_s, \text{se}(\hat{\beta}_s))$ for the EE model. These properties are potentially useful when collaborating among

groups, where sharing of raw data can be more difficult than sharing summary data.

2.4.2 Induced Single Study Bayes Factors

For the ES model, in the special case of one subgroup ($S = 1$), both the actual Bayes Factor and our approximations to it (2.23) reduce to results from previous work. More specifically, the approximation $\widehat{\text{BF}}^{\text{ES}}$ becomes exact in this case, and equal to the Bayes Factor derived by Servin and Stephens (2007), whereas ABF^{ES} is equal to the ABF in Wakefield (2009) (see also Johnson (2005) and Johnson (2008)).

2.4.3 Non-informative Subgroup Data

Suppose that in one of the subgroups, s , the sample genotypes concentrate in one of the three genotype categories. This situation might arise, for example, in cross-population genetic studies, where it is not uncommon to observe SNPs that are quite variable in some populations, but with little or no variation in other populations. Intuitively, in this case, subgroup s will contain little or no information for testing the global null hypothesis. Indeed, this will be correctly reflected in the standard error for the effect size, δ_s , being large, and this will in turn result in study s not contributing to the Bayes Factor. For example, it is easy to see that, for large δ_s (formally, in the limit $\delta_s \rightarrow \infty$), the ABF (2.23) is unaffected by the association data in study s (T_s); and a similar result holds for the exact Bayes Factor (appendix A) and for both EE and ES models. In other words, the Bayesian procedure correctly reflects the non-informativeness of the data from study s . Although this property seems very intuitive and one might expect every reasonable statistical procedure to possess it, many widely-used methods to combining results across studies do not have this property (e.g. Fisher’s combined probability test).

2.5 Extreme Models and Connections with Frequentist Tests

The proposed models are very flexible, covering a wide range of types and degrees of heterogeneity by setting different values for (ϕ, ω) (or k in the CEFN prior). In this section, we discuss the two extremes of no heterogeneity (“fixed effects”) and maximum heterogeneity, and establish connections between

Bayesian and frequentist testing approaches in these settings.

2.5.1 The Fixed Effects Model

A fixed effects model assumes genetic effects are homogeneous across subgroups. In the ES model this corresponds to setting $\phi = 0$, since this ensures $b_s = \bar{b}$ for all s . Similarly, in the EE model setting $\psi = 0$ ensures $\beta_s = \bar{\beta}$ for all s .

Since fixed effects analyses are widely used, we now briefly divert to connect our notation and models to existing methods. In the cases $\phi = 0$ and $\psi = 0$ the test statistics \mathcal{T}_{ES} and \mathcal{T}_{EE} defined above have particularly simple forms, being a weighted sum of the individual T_s statistics from each study (often referred to as a weighted sum of Z scores when the sample sizes are large). In particular, they are of the form

$$\mathcal{T} = \frac{\sum_s w_s T_s}{\sqrt{\sum_{s'} w_{s'}^2}} \quad (2.24)$$

where

1. For the ES model, $w_s = \text{se}(\hat{b}_s)^{-1} \approx \sqrt{2n_s f_s(1 - f_s)}$,
2. For the EE model, $w_s = \text{se}(\hat{\beta}_s)^{-1} \approx \hat{\sigma}_s^{-1} \sqrt{2n_s f_s(1 - f_s)}$,

and f_s denotes the allele frequency of the target SNP in subgroup s (with the approximations coming from an assumption of Hardy–Weinberg equilibrium in each subgroup). We note in passing that these representations give some insight into the main practical difference between the ES and EE models for testing: the EE model gives greater weight to studies with small residual error variance. Note also that the T_s values are the same for both EE and ES models, and independent of measurement scale, but σ_s depends on measurement scale, so \mathcal{T}_{ES} is robust to studies using different measurement scales (or different transformations of the phenotypes) but \mathcal{T}_{EE} is not. In addition, these representations clarify the connection between these statistics and the methods used in the meta-analysis software METAL (Willer et al. (2010)). Specifically, METAL implements tests using the weighted statistic (2.24) with

two different weighting schemes, one corresponding to the EE model weights above, and the other with the weights equal to $\sqrt{n_s}$. This latter scheme corresponds to the ES model only if f_s is equal across studies. (Where f_s varies across studies the weighting in the ES model seems, to us, preferable to the METAL scheme since, intuitively, studies with small $f_s(1 - f_s)$ provide less information.)

Returning now to the Bayes Factors, in the fixed effects cases, $\phi = 0$, the approximate Bayes Factors (2.23) simplifies to

$$\text{ABF}_{\text{fix}}^{\text{ES}}(\omega) := \text{ABF}^{\text{ES}}(\phi = 0, \omega) = \sqrt{\frac{\zeta^2}{\zeta^2 + \omega^2}} \cdot \exp\left(\frac{\mathcal{T}_{\text{ES}}^2}{2} \frac{\omega^2}{\zeta^2 + \omega^2}\right), \quad (2.25)$$

where

$$\zeta^2 = \frac{1}{\sum_s \delta_s^{-2}}. \quad (2.26)$$

And a similar expression holds for $\text{ABF}_{\text{fix}}^{\text{EE}}(w) := \text{ABF}^{\text{EE}}(\psi = 0, w)$.

Implicit p-value prior

We now answer the following question: under what prior assumptions will $\text{ABF}_{\text{fix}}^{\text{ES}}$ produce the same SNP rankings as the frequentist test statistic \mathcal{T}_{ES} ? Wakefield (2009) names this kind of prior the “implicit p -value prior”, as it effectively identifies the implicit prior assumptions that are made by the standard practice of ranking SNPs by their p -value computed from \mathcal{T}_{ES} .

Note that, although for a *given* SNP $\text{ABF}_{\text{fix}}^{\text{ES}}(\omega)$ is a monotone function of \mathcal{T}_{ES} , for a fixed value of ω the two statistics will not generally rank SNPs in the same way because ζ varies among SNPs. If however we allow ω to vary among SNPs in a particular way, then the two statistics will agree on their rankings, as summarized in the following proposition.

PROPOSITION 2. (Implicit p -value prior, fixed effects) *In the ES model, if the prior hyperparameter ω is allowed to vary among SNPs, with*

$$\omega_p = K\zeta_p, \quad (2.27)$$

where p indexes SNPs, and K is any positive constant, then $\text{ABF}_{\text{fix}}^{\text{ES}}$ and \mathcal{T}_{ES} will produce the same ranking of SNPs.

Proof. This follows directly from substituting (2.27) into (2.25). \square

Of course, a similar result holds for the EE model.

NOTE 3. Recall that ζ_p is the standard error of $\hat{\bar{b}}$ for SNP p , and so it is large if there is little information about \bar{b} at a particular SNP. In a meta-analysis setting ζ_p can vary across SNPs not only because allele frequencies may vary across SNPs, but also because data on some SNPs may be available in only a subset of studies. Recall also that large values of ω_p correspond to a prior assumption that the effect size \bar{b} at SNP p is likely to be large (in absolute value). Thus the implicit p -value prior effectively assumes that there are larger effects at SNPs with less information.

NOTE 4. In the idealized case where data on all SNPs are available on the same subgroups, and the subgroups also have similar allele frequencies at every SNP (as might happen if the subgroups come from a single random mating population) then the sample genotype variance of SNP p in subgroup s can be well approximated by $2n_s f_p(1 - f_p)$, where f_p is the population allele frequency of SNP p . (Note the slight abuse of notation, since we previously indexed f by subgroup, whereas here it is indexed by SNP.) Consequently, the implicit frequentist prior (2.27) can be written as

$$\omega_p = K \sqrt{\frac{1}{f_p(1 - f_p)} \frac{1}{\sum_s n_s}}, \quad (2.28)$$

which is effectively the same as the single subgroup case discussed by Wakefield (2009).

2.5.2 Maximum Heterogeneity Model

We now consider the other end of the spectrum: models with very high heterogeneity. Specifically, if we consider a class of ES models with some fixed prior expected marginal effect sizes ($\phi^2 + \omega^2$), then the model with $\omega = 0$ represents the maximum possible heterogeneity among all models in the class. In this case, the average effect \bar{b} is forced to be strictly 0, and the effects $b_s|\phi$ are independent, $N(0, \phi^2)$.

It can be shown from (A.13) and (A.28) that for both the EE and the ES models, the exact Bayes Factor under this particular setting, BF_{maxH} , is the product of the individual Bayes Factors, i.e.

$$\text{BF}_{\text{maxH}} = \prod_s \text{BF}_{\text{single}, s}, \quad (2.29)$$

where $\text{BF}_{\text{single}, s}$ is the exact Bayes Factor calculated using data only from subgroup s . This relationship also holds for the ABF, i.e.

$$\text{ABF}_{\text{maxH}}^{\text{ES}}(\phi) := \text{ABF}^{\text{ES}}(\phi, \omega = 0) = \prod_s \sqrt{\frac{\delta_s^2}{\delta_s^2 + \phi^2}} \exp\left(\frac{T_s^2}{2} \frac{\phi^2}{\delta_s^2 + \phi^2}\right). \quad (2.30)$$

To connect this with frequentist tests, consider the likelihood ratio test of the global null hypothesis $b_s = 0$ (for all s) vs the general alternative where the b_s values are unconstrained. The likelihood ratio test statistic for this test can be written as the product of likelihood ratio statistics in each subgroup:

$$\text{LR}_{\text{maxH}} = \prod_s \text{LR}_s, \quad (2.31)$$

where LR_s is the likelihood ratio test statistic for $b_s = 0$ vs the alternative that b_s is unconstrained. For reasonably large sample sizes, LR_s can be well approximated by

$$\lim_{n_s \rightarrow \infty} \text{LR}_s \approx \exp\left(-\frac{T_s^2}{2}\right), \quad (2.32)$$

and so

$$\text{LR}_{\text{maxH}} \approx \exp\left(-\frac{\sum_s T_s^2}{2}\right). \quad (2.33)$$

Thus the likelihood ratio test is approximately the same as a test based on $\sum_s T_s^2$, which (again assuming large sample sizes) is the sum of the squared Z values, and p -values can be obtained by noting that under the global null hypothesis this sum will be $\sim \chi_S^2$. This test is very similar to Fisher's approach to combining test statistics from multiple studies.

Now, again, we consider the following question: under what prior assumptions will $\text{ABF}_{\text{maxH}}^{\text{ES}}$ give the

same ranking of SNPs as $\sum_s T_s^2$? Under the ES model, we can see that in fact no single ϕ value will give this result. However, if we relax the ES model assumption by allowing ϕ to be subgroup specific, with

$$\phi_s^2 = K\delta_s^2, \tag{2.34}$$

where K is a constant for all subgroups and all SNPs tested, then it is easy to see that the resulting ABF yields the same ranking of SNP associations as $\sum_s T_s^2$.

NOTE 5. *Recalling that δ_s is the standard error for b_s , the implicit p-value prior (2.34) effectively assumes bigger effects in those subgroups with less information. There seems to be no good justification, in general, for this prior assumption.*

2.6 Extension to Case-Control Data

In situations when phenotypes are case/control status, we can replace the linear model (2.1) for each subgroup by a logistic regression model: for individual i in subgroup s , the phenotype-genotype association is modeled by

$$\log \frac{\Pr(y_{si} = 1|g_{si})}{\Pr(y_{si} = 0|g_{si})} = \mu_s + \beta_s g_{si}. \tag{2.35}$$

Furthermore, we use the same form for the priors on μ_s , β_s and $\bar{\beta}$ as described in the EE model for quantitative traits.

Computing Bayes Factors for this model is challenging because the marginal likelihood is analytically intractable. To ease the computation, we approximate the subgroup-level log-likelihood function $l(\beta_s, \mu_s)$ given by (2.35) using an asymptotic expansion around its maximum likelihood estimates. The resulting approximate Bayes Factor has the same functional decomposition form as in (2.23) only with t-statistics replaced by Wald statistics. We show the details of the derivation and the results in appendix B.

3 Applications and Illustrations

3.1 deCODE Recombination Study

The deCODE recombination study (Kong et al. (2008)) aimed to find genetic variants that explain genome-wide recombination rate variation. The study genotyped 1,887 males and 1,702 females from the Icelandic population and performed a genome-wide scan searching for association between SNP genotypes and the estimated genome-wide average recombination rate for each individual.

Prior to this study, it was already well known that male and female recombination rates differ considerably over moderate-sized genomic regions; the researchers therefore analyzed the data separately for males and females. They estimated the genetic effect sizes on the recombination phenotype assuming an additive model (2.1). For recombination rate in males, they found three highly correlated SNPs in a small region on chromosome 4p16.3 show strong association signals. Interestingly, these three SNPs also show strong associations in females, but with the estimated effect being in the opposite direction (i.e. the allele associated with lower recombination rate in males appears to be associated with higher recombination rate in females).

Here we use these data as a simple illustration of our Bayesian analysis tools, taking the summary-level data on three reported associations in Table 1 of Kong et al. (2008) to compute ABFs under different models of heterogeneity. In particular, we use their point estimates of effect sizes $\hat{\beta}_{\text{male}}$ and $\hat{\beta}_{\text{female}}$ and infer $\text{se}(\hat{\beta}_{\text{male}})$ and $\text{se}(\hat{\beta}_{\text{female}})$ from their corresponding reported p -values. With these summary data we are able to compute ABFs under the EE model. We treat males and females as two subgroups and consider 4 levels of expected marginal overall effect sizes with $\sqrt{\psi^2 + w^2} = 5, 10, 20, 40$ (the phenotype scale being centi-Morgans) and 5 levels of heterogeneity levels with $\psi^2/w^2 = 0, 0.5, 1, 2, \infty$. In total, we obtain a grid of 4×5 different (ψ, w) combinations and we treat every grid value as *a priori* equally likely when computing $\text{ABF}_{\text{av}}^{\text{EE}}$. (Of course it would be easy to consider a denser grid of values, and necessary if we wanted precise estimates of ψ and w , but this coarser grid suffices for our purposes here.)

The resulting Bayes Factors are shown in Table 1. Note that the meta-analysis Bayes Factor allowing for heterogeneity ($\text{ABF}_{\text{av}}^{\text{EE}}$) is many orders of magnitude larger than either of the subgroup-specific

Bayes Factors, which are themselves larger than the Bayes Factor under the fixed effects meta-analysis model ($\text{ABF}_{\text{fix}}^{\text{EE}}$). This illustrates two simple but important points. First, that a meta-analysis can yield considerably stronger signal than subgroup-specific analyses, and second that in this case a standard fixed-effects meta-analysis would be ineffective at identifying the association signal. Of course, in general, the “right” level of heterogeneity is unknown; $\text{ABF}_{\text{av}}^{\text{EE}}$ deals with this problem by averaging over different levels of heterogeneity (including the fixed effects, or no heterogeneity, case). In general we view this ability to average over unknown quantities as an attractive feature of the Bayesian approach, although as we shall see in the next application it can be helpful to examine the components of this average separately.

SNP	Male		Female		Meta Bayes Factors	
	Effect (p-value)	$\text{ABF}_{\text{single,male}}^{\text{EE}}$	Effect (p-value)	$\text{ABF}_{\text{single,female}}^{\text{EE}}$	$\text{ABF}_{\text{fix}}^{\text{EE}}$	$\text{ABF}_{\text{av}}^{\text{EE}}$
rs3796619	−67.9 (1.1×10^{-14})	$10^{11.12}$	67.6 (7.9×10^{-6})	$10^{2.81}$	$10^{3.07}$	$10^{13.91}$
rs1670533	−66.1 (1.8×10^{-11})	$10^{8.06}$	92.8 (4.1×10^{-8})	$10^{4.55}$	$10^{1.10}$	$10^{12.58}$
rs2045065	−66.2 (1.6×10^{-11})	$10^{8.11}$	92.2 (6.0×10^{-8})	$10^{4.40}$	$10^{1.18}$	$10^{12.49}$

Table 1: Bayesian meta-analysis of genetic associations with recombination rate. The SNPs and their estimated effect sizes and p-values are directly taken from Kong et al. (2008) Table 1. We compute approximate Bayes Factors under the EE model using only those reported summary statistics. $\text{ABF}_{\text{single,male}}^{\text{EE}}$ and $\text{ABF}_{\text{single,female}}^{\text{EE}}$ are approximate Bayes Factors computed using only male subgroup and female subgroup data respectively.

Having established that an association exists, one might turn to assessing the extent of the heterogeneity among subgroups. We note that, although the p -values of the three SNPs indicate the genetic effects in males and females are separately both highly significant, and the effect size estimates have opposite signs, the p -values themselves do not directly assess evidence for heterogeneity. In contrast, the Bayes Factors in Table 1 can be compared with one another to help directly assess the evidence for heterogeneity. Indeed, the fact that $\text{ABF}_{\text{av}}^{\text{EE}}$ is substantially larger than $\text{ABF}_{\text{fix}}^{\text{EE}}$ immediately indicates that the data are inconsistent with the fixed effects assumption that effects are the same in both subgroups. More detailed investigation can be performed by comparing the Bayes Factors computed under different levels of the heterogeneity parameter (ψ^2/w^2); in this case the data are consistent with infinite values for this parameter (i.e. with the “maximum heterogeneity model”, $w^2 = 0$).

On the face of it, the analysis above appears to support the conclusion that one of these SNPs may have causal effects of different signs in males and females. However, this may not be the case. In particular, a recent association analysis of recombination rates in other samples suggested that this genetic region may actually contain more than one genetic variant affecting recombination rates, some acting in males and others in females, rather than a single genetic variant with antagonistic effects in the two groups (Fledel-Alon et al. (2011)). A more detailed assessment of this would require association analysis of multiple SNPs simultaneously, rather than the single SNP analyses we consider here.

3.2 Global Lipids Study

The Global Lipids consortium (Teslovich et al. (2010)) conducted a large scale meta-analysis of genome-wide genetic association studies of blood lipids phenotypes. In this study, more than 100,000 individuals of European ancestry were amassed through 46 separate studies (grouped into 25 studies in their final analysis). For each individual, measures of total cholesterol (TC), low-density lipoprotein cholesterol (LDL-C), high-density lipoprotein cholesterol (HDL-C) and triglycerides (TG) were obtained. Large-scale (whole genome) genotyping of genetic variants (SNPs) was performed and missing genotypes were imputed: in total, about 2.7 million common SNPs were included in the final association analysis. In each individual study, all four phenotypes were independently quantile normal transformed; single SNP association tests were performed for all SNPs and all phenotypes using the linear model (2.1) and the estimated effect sizes and their standard errors were computed. The meta-analysis combined these summary-level data from each individual study using the software METAL (Willer et al. (2010)), using the weighted statistic (2.24) with weights $w_s = \sqrt{n_s}$, which, as noted above, can be viewed as an approximation to a fixed effects analysis under the ES model if the SNP genotype frequencies are similar across studies. Teslovich et al. (2010) reported 168 SNP-phenotype associations exceeding their “genome-wide significant” threshold (fixed effects p -value $< 5 \times 10^{-8}$), and identified 95 genes, with 59 showing genome-wide significant association for the first time.

Given that this meta-analysis involves 25 different subgroups, with widely differing enrollment criteria (e.g. different disease cohorts), one might expect that heterogeneity of effects could be an issue. Indeed, many of the reported associated loci had effect size estimates of different signs in different studies.

Furthermore, one might worry that the fixed effects analysis used in the original analysis could have missed significant associations with moderate heterogeneity. To assess these issues we reanalyzed the data, using our Bayesian tools that allow for different degrees of heterogeneity.

To perform these analyses we were able to obtain access to summary data from each study, in the form of an estimated effect size (computed from the quantile-transformed phenotype data) and its standard error for each SNP in each study. These summary data allow us to perform analyses under the EE model, rather than the ES model effectively used in the original analysis. We present results for ABFs computed under three different types of EE model that assume increasing amounts of heterogeneity: the fixed effects model, the “limited heterogeneity” CEFN model, and the maximum heterogeneity model. In each case we assumed a discrete uniform prior on the overall genetic effect size (i.e., $(k^2 + 1)w^2$ in the CEFN model and $w^2 + \psi^2$ in the other models) on the set $\{0.1^2, 0.2^2, 0.4^2, 0.6^2, 0.8^2\}$. For the fixed effects model, $\psi = 0$; for the maximum heterogeneity model, $w = 0$; and for the CEFN model we set $k = 0.326$ which gives a prior probability of 1/1,000 that the genetic effect in each study has an opposite sign to the expected average effect $\bar{\beta}$.

Although a fully automated Bayesian analysis would most naturally average over the three different models of heterogeneity we considered, in practice, because of their different sensitivity to particular data features (see below) we found it helpful to examine them separately in this application.

As an initial check on data handling we verified that our Bayesian fixed effects analysis produced similar results to the original fixed effects analysis. As expected, we found that $\text{ABF}_{\text{fix}}^{\text{EE}}$ ranked SNP associations very similarly to the original reported (ES model) p -values. However, there were a few notable exceptions. In particular, a few SNPs showed much stronger association signal in $\text{ABF}_{\text{fix}}^{\text{EE}}$ than the original analysis. Further investigation suggested that these results likely reflected the EE analysis being less robust to certain issues than the original ES analysis (rather than, say, due to the difference between the prior distributions we assume in our Bayesian analysis and the implicit prior distributions assumed by the original p -value analysis). For example, SNP rs17061870 with LDL phenotype, had $\text{ABF}_{\text{fix}}^{\text{EE}} > 10^{17}$ ($\mathcal{T}_{\text{EE}}^2 = 86.46$ with a corresponding p -value $= 1.4 \times 10^{-20}$) compared with an original reported p -value $= 0.028$ ($\mathcal{T}_{\text{ES}}^2 = 4.85$). But examination of the study-specific data for this SNP showed a suspicious pattern: the p -value was 2×10^{-31} in one study (the Family Heart Study, FHS), but no

smaller than 0.1 in the 5 other studies for which this SNP had genotype data available. Furthermore, the very small p -value in FHS was driven primarily by a very small, probably erroneous, estimate of the residual error in that study (the sample size of this particular study is not large) which under the EE model results in a very high weight on that study, but in the ES model does not. We emphasize that we performed the EE analysis here because it was what we were able to do easily with the available summary data, rather than because we prefer it.

Next we assessed the evidence for heterogeneity of effects in the 168 association signals reported by the original analysis. We did this by comparing the support for the limited heterogeneity model ($ABF_{\text{cefn}}^{\text{EE}}$) with the support for the no heterogeneity model ($ABF_{\text{fix}}^{\text{EE}}$). The majority of phenotype-SNP pairs (111/168) showed stronger support for the no heterogeneity model, although a few showed overwhelming support for the limited heterogeneity model (Figures 1, 2). In light of the discussion of the recombination example above, it is important to note that even these apparently overwhelming signals for heterogeneity may actually reflect other factors (e.g. multiple SNPs in a gene affecting phenotype, or different measurement protocols in different studies) rather than genuine biological interactions with study group.

Of course, the fact that the 168 associations reported in the original analysis generally showed little evidence for heterogeneity could be explained by the fact that these associations were identified by a fixed-effects analysis that assumed no heterogeneity. Therefore, we finally assessed whether the original fixed effects analysis likely missed some associations showing heterogeneity across studies. To do this, we performed a genome-wide analysis for each phenotype, in each case excluding all SNPs within 1Mb of any SNPs originally reported as being associated with that phenotype, and searching for SNPs that showed strong evidence for association under one of the heterogeneity models ($ABF_{\text{cefn}}^{\text{EE}}$ or $ABF_{\text{maxH}}^{\text{EE}} \geq 10^6$) but not under the fixed-effects model ($ABF_{\text{fix}}^{\text{EE}} < 10^6$). This threshold (10^6) corresponds very roughly to, and is perhaps slightly more conservative than, the threshold effectively used in the initial analysis. (We used the same threshold for all three models for simplicity, but in this setting it would probably be better to use a variable threshold; for example, it would probably be appropriate to use a more stringent threshold for $ABF_{\text{maxH}}^{\text{EE}}$ than the other models, reflecting the fact that strong heterogeneity is uncommon in this context.)

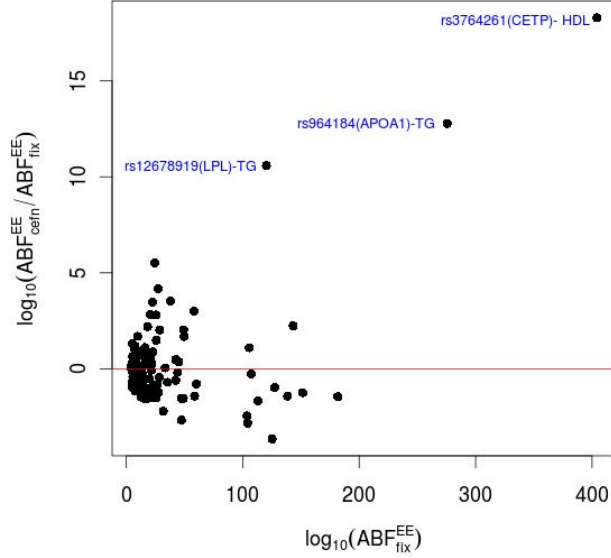


Figure 1: Assessment of evidence for heterogeneity in 168 reported phenotype-SNP association signals from Teslovich et al. (2010). Large values on the y axis indicate stronger support for the model allowing for limited heterogeneity ($\text{ABF}_{\text{cefn}}^{\text{EE}}$) compared with the model with no heterogeneity ($\text{ABF}_{\text{fix}}^{\text{EE}}$). The three highlighted points correspond to associations with overwhelming evidence for the limited heterogeneity model, $\text{ABF}_{\text{cefn}}^{\text{EE}} / \text{ABF}_{\text{fix}}^{\text{EE}} > 10^{10}$; forest plots for these three associations are shown in Figure 2

Overall we found 42 SNP-phenotype associations satisfying this criteria (counting multiple signals due to SNPs in LD with one another as a single association), representing associations potentially missed by the original analysis. However, detailed investigation of these suggested to us that many of them were likely not to be genuine associations. For example, 36 of these associations showed strong signals in $\text{ABF}_{\text{maxH}}^{\text{EE}}$, but not in the other models; but all these were driven by apparently strong associations in the FHS study alone, which themselves seemed likely to be due to data processing errors (e.g. for these SNPs the p -values in the FHS for quantile transformed phenotypes were many orders of magnitude smaller than for the original phenotypes). We then dropped the FHS data and re-performed this analysis. All 6 remaining signals from the previous analysis still satisfy the criteria (Table 2). Of those, the first two listed seem almost certain to be genuine: the genes ABCA1 and KLF14 are reported in Teslovich et al. (2010) as associated with other lipid phenotypes (ABCA1 with HDL and TC; KLF14 with HDL), but not with the phenotypes we listed in Table 2, and both reflect associations that just missed being significant in the original fixed effects analysis. The next two map (in the human genome build 36.3 that

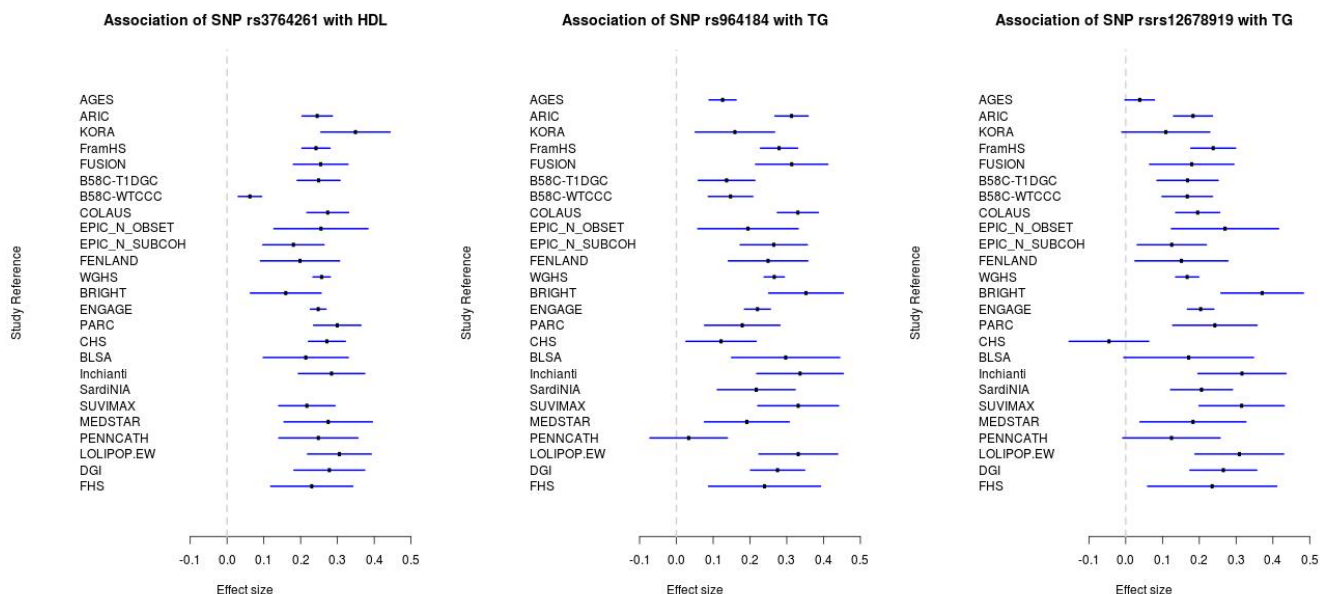


Figure 2: Forest plots for the three highlighted association signals in Figure 1, which showed overwhelming evidence for heterogeneity of (apparent) effects.

we used) approximately 6Mb apart on chromosome 11; however examination of the raw data showed these SNPs to have similar effect estimates across studies, suggesting that they are correlated with one another. Usually SNPs 6Mb apart would be uncorrelated with one another, but this could be explained by mapping errors, and indeed in more recent genome build the SNPs are closer together (though still 2.5Mb apart). Further, both SNPs map to within a few Mb of the gene *LRP4*, which was reported as strongly associated with HDL in the original analysis. It seems quite possible that the signals at these SNPs are simply reflecting correlation of these SNPs with this reported association and not a novel association. Finally, the last two associations are again driven by apparent anomalies in a single study (this time B58C-WTCCC).

In summary, we find that the original fixed effects analysis in Teslovich et al. (2010) was highly effective at identifying the main association signals in these data. Nonetheless, some of the (presumably real) associations identified by the original study do show a substantially stronger signal in analyses that allow for heterogeneity (Figure 1), and it remains possible that in other data sets meta-analyses allowing for heterogeneity could cause findings that are “borderline significant” under a fixed effects model to be promoted to a level of convincing evidence. On the other hand, our results also provide a cautionary

Phenotype	SNP	Gene Region	$\log_{10}(\text{BF}_{\text{fix}}^{\text{EE}})$	$\log_{10}(\text{BF}_{\text{maxH}}^{\text{EE}})$	$\log_{10}(\text{BF}_{\text{cefn}}^{\text{EE}})$
LDL	rs1800978	5'UTR of ABCA1	5.2	3.4	6.0
TG	rs1562398	Flanking KLF14	5.3	-0.2	6.5
HDL	rs11229165	Flanking OR4A16	4.6	4.9	6.4
HDL	rs7108164	Flanking OR4A42P	4.2	4.9	6.3
HDL	rs11984900	N.A.	-1.1	16.6	6.2
HDL	rs6995137	Flanking SFRP1	-0.4	6.9	4.8

Table 2: Association signals that show strong association under the models allowing for heterogeneity ($\text{ABF}_{\text{cefn}}^{\text{EE}}$ or $\text{ABF}_{\text{maxH}}^{\text{EE}} \geq 10^6$) but less strong under a model with no heterogeneity ($\text{ABF}_{\text{fix}}^{\text{EE}} < 10^6$). It seems likely that only the first two of these reflect genuine associations missed by the original analysis (see text for discussion).

tale: in the context of meta-analysis of genetic association studies, when associations appear only under models allowing for strong heterogeneity, and not under fixed effects models, the reasons for the discrepancy must be examined carefully, and the results interpreted critically. Indeed, it seems that searching for SNPs showing strong apparent heterogeneity of effects provides an effective way to identify data processing errors that may otherwise lurk undetected!

3.3 Population eQTL Study

In this section, we apply the proposed Bayesian methods to map expression quantitative trait loci (eQTLs) using multi-population data. An eQTL is a genetic variant (here we only focus on SNPs) that is associated with gene expression phenotype.

We consider gene expression measurements from Stranger et al. (2007), obtained using the Illumina Sentrix Human-6 Expression BeadChip, on lymphoblastoid cell lines derived from 141 unrelated individuals from the International Hapmap project (The International HapMap Consortium (2005)). These individuals were sampled from three major population groups (41 Europeans (CEU), 59 Asians (ASN) and 41 Africans (YRI)) and were fully sequenced in the pilot project of the 1000 Genomes project (Durbin et al. (2010)). We focus on the 8,427 distinct autosomal genes that were confirmed to be expressed in the same African samples by an independent RNA-seq experiment (Pickrell et al. (2010)). We used the SNP genotype data on 14.4 million SNPs from the final release (March, 2010) of the pilot SNP calls

from the 1000 genomes project, with no additional allele frequency filtering applied. In addition to the original normalization applied to the expression data in Stranger et al. (2007), we perform quantile normal transformations to each selected gene, separately within each population group, to reduce the influence of outliers or other deviations from normality.

Previous studies have shown that most eQTLs are located near to the gene whose expression they influence (so-called “*cis*-eQTLs”). Therefore for each gene we restrict our association analysis to the “*cis* SNPs” which lie within the region 500kb upstream of the transcription start site and 500 kb downstream of the transcription end site. Previous analyses of these kinds of data have also either ignored potential heterogeneity of effects across populations, or have examined heterogeneity by analyzing each population separately and then looking for differences between the results in each population (e.g. Stranger et al. (2007)). However, analyzing populations separately in this way has several disadvantages: for example, it has less power to identify eQTLs than a joint analysis (at least, those eQTLs exhibiting little heterogeneity), and if an eQTL is identified in one population but not another then it is unclear whether this likely represented genuine heterogeneity or lack of power. Our Bayesian tools help overcome these problems: we can analyze all the populations together, allowing for potential heterogeneity, and also assess how strong the evidence is for heterogeneity in the eQTLs identified.

To perform our analyses, we group the samples by their population of origin to form three subgroups and apply the proposed Bayesian methods. We use the ES model with a grid of (ϕ, ω) values. Specifically, we consider five levels of $\sqrt{\phi^2 + \omega^2}$ values: 0.1, 0.2, 0.4, 0.8, 1.6, and seven degrees of heterogeneity characterized by ϕ^2/ω^2 values: 0, 1/4, 1/2, 1, 2, 4, ∞ . Further, we assign these 35 grid values equal prior weight. These prior distributions are broad, covering a wide range of possible effect sizes and levels of heterogeneity. For each gene, we compute $\widehat{\text{BF}}_{\text{av}}^{\text{ES}}$ for each *cis* SNP and identify the highest ranked (“top”) SNP. For each top SNP we assess the evidence for heterogeneity by examining how the Bayes Factors vary with the heterogeneity ϕ^2/ω^2 .

We begin by quantifying the fact that mapping eQTLs in all groups simultaneously produces stronger association signals than considering the groups separately. Specifically, for the most strongly associated SNP in each gene, we compare the overall Bayes Factor for association averaging over many different possible levels of heterogeneity ($\widehat{\text{BF}}_{\text{av}}^{\text{ES}}$) with the largest of the Bayes Factors obtained from

considering each population separately ($\text{BF}_{\text{max_single}}^{\text{ES}}$). Among 5,691 genes in which top associated SNPs are polymorphic in at least 2 populations, 4,470 of them (79%) show stronger signal when mapping together ($\widehat{\text{BF}}_{\text{av}}^{\text{ES}} > \text{BF}_{\text{max_single}}^{\text{ES}}$) and 2,214 of them (39%) show considerably stronger signal ($\widehat{\text{BF}}_{\text{av}}^{\text{ES}} / \text{BF}_{\text{max_single}}^{\text{ES}} > 10$).

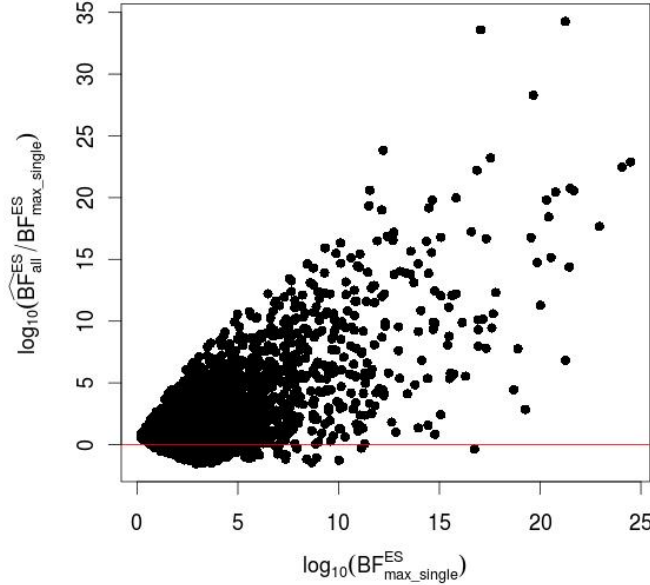


Figure 3: Comparison of strength of association signal from mapping eQTLs separately in each population ($\text{BF}_{\text{max_single}}^{\text{ES}}$), and simultaneously in all populations, allowing for heterogeneity ($\widehat{\text{BF}}_{\text{av}}^{\text{ES}}$). A clear majority of associations (79%) show stronger association in the joint analysis, often much stronger.

Examining the posterior distributions of the heterogeneity parameter ϕ^2/ω^2 for each SNP suggested that most eQTLs showed little heterogeneity across populations (results not shown). Here we focus on the relatively few SNPs that showed evidence for strong heterogeneity; specifically we focus on SNPs for which the Bayes Factor for the maximum heterogeneity model ($\phi^2/\omega^2 = \infty$) versus the no-heterogeneity model ($\phi^2/\omega^2 = 0$) is largest. Examining the value of this ratio across all top SNPs (Figure 4) we identified 10 SNPs with very strong evidence in favor of the maximum heterogeneity model ($\widehat{\text{BF}}_{\text{maxH}}^{\text{ES}} / \widehat{\text{BF}}_{\text{fix}}^{\text{ES}} > 10^5$). These are listed in Table 3, and four are illustrated in more detail in Figure 5.

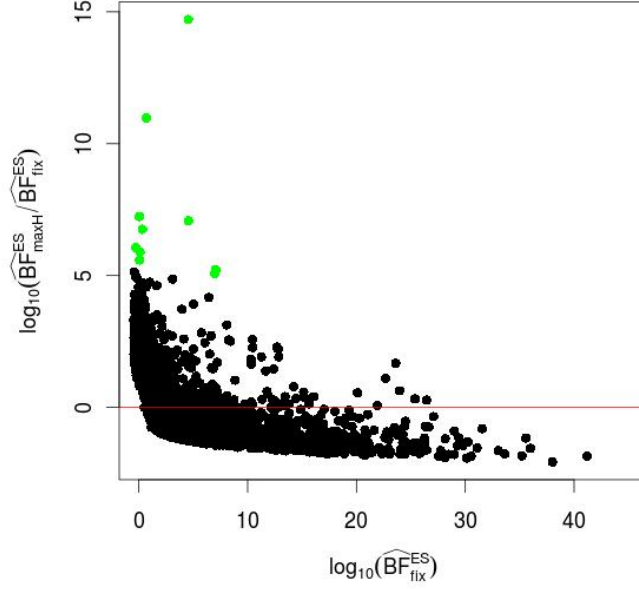


Figure 4: Comparisons of maximum heterogeneity models and fixed effects models for all top ranked SNPs. The 10/8,427 SNPs with strong evidence for the heterogeneity model ($\widehat{\text{BF}}_{\text{maxH}}^{\text{ES}} / \widehat{\text{BF}}_{\text{fix}}^{\text{ES}} > 10^5$ and $\widehat{\text{BF}}_{\text{av}}^{\text{ES}} > 10^5$) are highlighted in green, and listed in Table 3.

SNP	Gene	$\log_{10}(\widehat{\text{BF}}_{\text{fix}}^{\text{ES}})$	$\log_{10}(\widehat{\text{BF}}_{\text{maxH}}^{\text{ES}})$	$\log_{10}(\widehat{\text{BF}}_{\text{av}}^{\text{ES}})$
rs9595893	RP11-298P3.4	4.5	19.3	19.0
rs1037495	SOS1	7.1	12.3	12.1
rs3180068	PAQR8	7.0	12.0	12.0
rs380359	PLA2G4C	4.7	11.6	11.5
rs6008545	TTC38	0.7	11.7	11.3
rs1313996	C6orf48	0.0	7.3	7.0
rs11070253	BUB1B	0.3	7.1	6.8
rs4072597	HEATR2	0.1	6.0	5.7
rs7581360	POLR2D	-0.3	5.8	5.5
rs437380	C17orf48	0.1	5.7	5.4

Table 3: Details of the 10 eQTL SNPs showing strongest evidence for the maximum heterogeneity model.

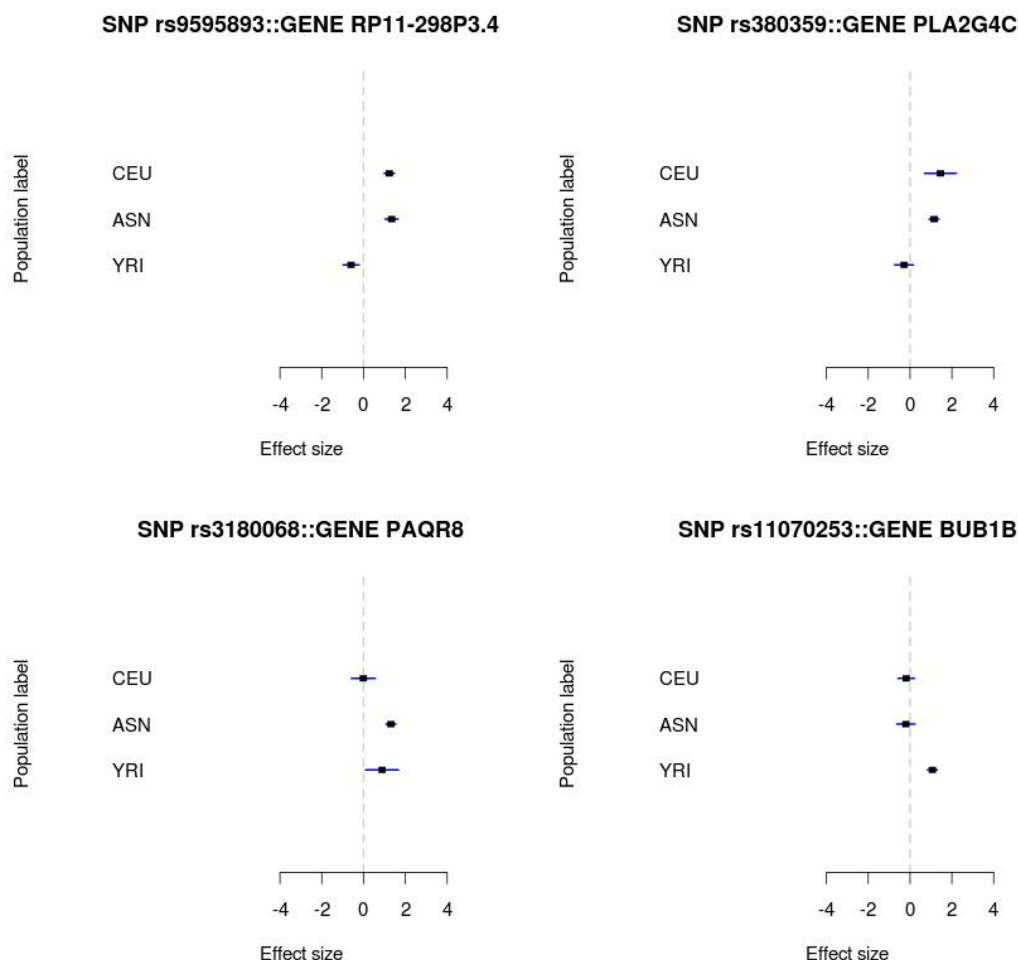


Figure 5: Forest plots for four of the 10 eQTL SNPs showing strongest evidence for the maximum heterogeneity model. Each panel shows the estimated effect and its 95% confidence interval in each population (obtained from a standard linear regression applied separately in each population).

3.3.1 Investigating Potential Population-specific eQTLs

Examining Figure 5 suggests that at least some of the eQTLs with evidence for heterogeneity may be consistent with a particular type of heterogeneity: having an effect in only some populations, with no effect in others. That is they might be “population-specific”.¹ For example, the forest plot for the eQTL in BUB1B suggests this eQTL may be specific to the YRI sample. This kind of heterogeneity could

¹It could be objected that the notion of a “population-specific” eQTL is too simplistic, and that apparent absence of effects in some populations more likely reflects very small, non-zero, effects. While sympathetic to this argument, we also find the simplicity has a certain appeal, and we view such models as potentially useful nonetheless.

occur if, for example, an eQTL affected binding of a particular regulatory element that is active in some populations and not others. Such differences in activity could, for example, happen as a consequence of natural selection (Kudaravalli et al. (2009)).

Motivated by this, we here describe how we can use the models described above to explicitly allow for “subgroup-specific” effects. Let C denote a binary string of indicators for whether an eQTL is active (i.e. has non-zero genetic effect) in each population. For example $C = (110)$ would indicate that the eQTL is active only in the first two populations. For the three populations in our data, C has 2^3 possible values, which we refer to as “configurations”. To evaluate the support for a particular configuration, we compute a Bayes Factor that contrasts the marginal likelihood of the configuration of interest to the null model $C = (000)$. For example, for $C = (110)$,

$$\begin{aligned} \text{BF}_{C=(110)} &= \frac{P(\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3 | \mathbf{g}_1, \mathbf{g}_2, \mathbf{g}_3, C = (110))}{P(\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3 | \mathbf{g}_1, \mathbf{g}_2, \mathbf{g}_3, C = (000))} \\ &= \frac{P(\mathbf{y}_1, \mathbf{y}_2 | \mathbf{g}_1, \mathbf{g}_2)}{P(\mathbf{y}_1, \mathbf{y}_2 | H_0)}. \end{aligned} \tag{3.1}$$

(The simplification is due to the assumption that the vectors of residual errors in (2.1) are independent across populations.) Note that this BF depends only on the data from the subgroups where the eQTL is assumed active.

Here, for the populations in which the eQTL is active, we use the ES model and the CEFN prior (with $k = 0.314$), so that we are allowing for only “limited” heterogeneity of effects in populations where the eQTL is active (which might be expected *a priori*, and is also consistent with the examples in Figure 5).

To illustrate this approach we apply it to SNP rs11070253 in the gene BUB1B. The BFs for all eight configurations, Table 4, show strong support for the configuration where this eQTL is active only in the Yoruban population.

It should be evident that a lot more work is required to turn these kinds of observations into solid biological insights. Indeed, without additional experimental confirmation we regard it as unclear whether any of the eQTLs we highlight here represent genuine interactions. Nonetheless, we believe this example

CEU	ASN	YRI	\log_{10} BF
0	0	0	0.0
1	0	0	-0.3
0	1	0	-0.3
0	0	1	8.1
1	1	0	-0.3
1	0	1	5.8
0	1	1	5.9
1	1	1	3.8

Table 4: Evaluation of population specificity for SNP rs11070253 and gene BUB1B. The \log_{10} Bayes Factors for all possible activity configurations (the left column) are shown.

is helpful both for illustrating the kinds of analysis that are possible with the tools we have developed, and also for highlighting putative interacting eQTLs that may be interesting for further study.

4 Discussion

The main contribution of this paper is to provide a flexible toolbox of methods for Bayesian association analyses involving potentially heterogeneous subgroups. We focus particularly on the genetic association context, which is characterized by the fact that the first goal is generally to reject the “global null”; that is, to identify genetic variants that are associated with phenotype of interest in some subgroup. Our examples illustrate how the tools we present can be used to i) identify associations *allowing for* different amounts and types of heterogeneity; and ii) investigate the relative strength of the evidence for different amounts and types of heterogeneity. The tools are sufficiently flexible to tackle a wide range of applications, including both applications involving limited heterogeneity as one might expect in typical meta-analyses, and the more extreme heterogeneity that one might encounter in the context of a Gene-Environment interaction, or subgroup-specific effects. We have developed different computationally efficient approximations to obtain numerical results, which is critical for handling of large scale genetic association data. In addition, we have highlighted connections between Bayes Factors and standard frequentist test statistics in this context (e.g. Propositions 1 and 2).

Many of the models and priors considered here have close connections with those employed in the

context of meta-analysis. In particular, they are very similar to the mixed effect meta-analysis models in standard frequentist approaches for studying of quantitative phenotypes, where the subgroup-specific intercept terms μ_s in (2.1) are regarded as fixed effects terms and genetic effect β_s (or b_s) are regarded as random effect terms.

Our models are also connected with, but differ in an important way from, those most frequently used in studies of gene-environment (G×E) interactions. The typical model used in this context is a linear model with marginal effect of subgroups and an gene-subgroup interaction term included, i.e.

$$y_i = \mu + \beta_e s_i + \beta_g g_i + \beta_{[g:e]} s_i g_i + e_i, \quad e_i \sim N(0, \sigma^2), \quad (4.1)$$

where s_i is a dummy variable denoting the subgroup membership of individual i , and $\beta_{[g:e]}$ is the coefficient of the subgroup-genotype interaction term and often of interest. This model is quite similar to (2.1). By re-arranging and grouping the terms, the linear model can be written as

$$y_i = (\mu + \beta_e s_i) + (\beta_g + \beta_{[g:e]} s_i) g_i + e_i, \quad e_i \sim N(0, \sigma^2). \quad (4.2)$$

Essentially, each subgroup is described with its own intercept, $\mu + \beta_e s_i$, and its own genetic effect, $\beta_g + \beta_{[g:e]} s_i$. (Note, if a marginal effect of subgroup is not included, the model is making a much stronger assumption on equal intercepts for different subgroups, which can be dangerous in practice and may lead to Simpson’s paradox (Bravata and Olkin (2001))). Nevertheless, the interaction model still makes stronger assumptions by assuming that the error variances across subgroups are the same. In comparison, our models allow this quantity to vary across in subgroups. This robust assumption is highly desirable in genetic association context, because, most likely, neither model (2.1) nor model (4.1) captures all factors affecting the phenotype of interest, and confounding factors almost certainly exist. In our meta-analysis model, the effect of the unaccounted confounding factors are “absorbed” by both the intercept and error variance terms, whereas the interaction model, lacking flexible error variance terms, does not possess this property.

One important issue that we have largely ignored here is the question of how to weigh evidence of heterogeneity in the data (e.g. large Bayes Factors for high heterogeneity models) against an *a priori*

belief that, in general, strong heterogeneity might be rare. In principle this is straightforward: given a prior distribution on different types of heterogeneity, it is trivial to use the Bayes Factors to compute posterior distributions. However, there remains an issue of choice of appropriate priors (an issue that also arises in a disguised form in frequentist approaches, for example in selecting appropriate p -value thresholds when testing for heterogeneity). Here we have often used (discrete) uniform distributions for convenience. In general, one might certainly want to change this, and appropriate priors may be context-dependent. For example, in a meta-analysis context, one might put more weight on models allowing for limited heterogeneity (the CEFN prior), whereas in an gene-environment interaction context one might be willing to allow more heterogeneity.

An attractive alternative to specifying context-specific priors on different levels of heterogeneity would be to attempt to “learn” how common are different levels of heterogeneity from the available data. This could be done by incorporating the Bayes Factors we have developed here into a hierarchical model, whereby, given appropriate data, the relative frequency of different types of heterogeneity could be estimated from the data. This approach might be particularly helpful for gene expression studies, where the large number of simultaneously measured phenotypes makes it plausible that the data could be quite informative for the relative frequency of different types of heterogeneity, and we see this as an important area for future work.

5 Acknowledgement

This work was supported by NIH grants HG02585 to M.S. and MH090951-02 (PI Jonathan Pritchard). We thank Yongtao Guan, Tanya Teslovich, Daniel Gaffney, Michael Stein and Peter McCullagh for helpful discussions. We thank the Global Lipids Consortium for access to their summary data.

A Computing Bayes Factors

In this section, we show the detailed calculations of various Bayes Factors.

A.1 Computation in the ES Model

A particular ES model, describing an alternative hypothesis H_a , is fully specified by setting values for (ϕ, ω) and hyper-parameters $(v_1, \dots, v_S, l_1, m_1, \dots, l_S, m_S)$. Under the contrasting null model H_0 , we set $\phi = \omega = 0$ while keeping all other hyper-parameters the same.

Let $\boldsymbol{\beta}_s = (\mu_s, \beta_s)$, $\tau_s = \sigma_s^{-2}$ and $\boldsymbol{\theta} = (\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_S, \tau_1, \dots, \tau_S, \bar{b})$, the marginal likelihood under model H_a can be written as

$$\begin{aligned} P(\mathbf{Y}|\mathbf{G}, H_a) &= \int P(\mathbf{Y}|\mathbf{G}, \boldsymbol{\theta}, H_a) p(\boldsymbol{\theta}|H_a) d\boldsymbol{\theta} \\ &= \int \left(\prod_s P(\mathbf{y}_s|\mathbf{g}_s, \boldsymbol{\beta}_s, \tau_s) \prod_s P(\boldsymbol{\beta}_s|\tau_s, \bar{b}, H_a) \prod_s P(\tau_s|H_a) P(\bar{b}|H_a) \right) d\boldsymbol{\beta}_1 \cdots d\boldsymbol{\beta}_S d\tau_1 \cdots d\tau_S d\bar{b} \quad (\text{A.1}) \\ &= \int \left(\int \left(\prod_s \int P(\mathbf{y}_s|\mathbf{g}_s, \boldsymbol{\beta}_s, \tau_s) P(\boldsymbol{\beta}_s|\tau_s, \bar{b}, H_a) d\boldsymbol{\beta}_s \right) p(\bar{b}|H_a) d\bar{b} \right) \prod_s P(\tau_s|H_a) d\tau_1 \cdots d\tau_S \end{aligned}$$

Let $\mathbf{X}_s = (\mathbf{1} \ \mathbf{g}_s)$ denote the design matrix of regression model (2.1) for subgroup s , it follows that

$$\begin{aligned} P(\mathbf{y}_s|\mathbf{g}_s, \boldsymbol{\beta}_s, \tau_s) &= \left(\frac{2\pi}{\tau_s} \right)^{-n_s/2} \exp \left(-\frac{\tau_s}{2} (\mathbf{y}_s - \mathbf{X}_s \boldsymbol{\beta}_s)' (\mathbf{y}_s - \mathbf{X}_s \boldsymbol{\beta}_s) \right) \\ &= \left(\frac{2\pi}{\tau_s} \right)^{-n_s/2} \exp \left(-\frac{1}{2} (\tilde{\mathbf{y}}_s - \mathbf{X}_s \mathbf{b}_s)' (\tilde{\mathbf{y}}_s - \mathbf{X}_s \mathbf{b}_s) \right), \end{aligned} \quad (\text{A.2})$$

where $\tilde{\mathbf{y}}_s = \sqrt{\tau_s} \mathbf{y}_s$ and $\mathbf{b}_s = \sqrt{\tau_s} \boldsymbol{\beta}_s = (\sqrt{\tau_s} \mu_s, b_s)$. We further denote

$$\bar{\mathbf{b}} = \begin{pmatrix} 0 \\ \bar{b} \end{pmatrix} \quad \text{and} \quad \Phi_s = \begin{pmatrix} v_s^2 & 0 \\ 0 & \phi^2 \end{pmatrix}, \quad (\text{A.3})$$

and write prior distribution $P(\mathbf{b}_s|\bar{\mathbf{b}}, H_a)$ in following matrix form,

$$\mathbf{b}_s|\bar{\mathbf{b}}, H_a \sim N(\bar{\mathbf{b}}, \Phi_s). \quad (\text{A.4})$$

We compute the marginal likelihood by sequentially evaluating the following integrals,

$$\begin{aligned}
F_{H_a,s} &= \int P(\mathbf{y}_s | \mathbf{X}_s, \mathbf{b}_s, \tau_s) P(\mathbf{b}_s | \bar{\mathbf{b}}, H_a) d\mathbf{b}_s \\
&= \left(\frac{2\pi}{\tau_s}\right)^{-n_s/2} |\Phi_s|^{-\frac{1}{2}} \cdot |\mathbf{X}_s' \mathbf{X}_s + \Phi_s^{-1}|^{-\frac{1}{2}} \\
&\quad \cdot \exp\left(-\frac{1}{2} \left(\tilde{\mathbf{y}}_s' \tilde{\mathbf{y}}_s - (\mathbf{X}_s' \tilde{\mathbf{y}}_s + \Phi_s^{-1} \bar{\mathbf{b}})' (\mathbf{X}_s' \mathbf{X}_s + \Phi_s^{-1})^{-1} (\mathbf{X}_s' \tilde{\mathbf{y}}_s + \Phi_s^{-1} \bar{\mathbf{b}}) + \bar{\mathbf{b}}' \Phi_s^{-1} \bar{\mathbf{b}}\right)\right).
\end{aligned} \tag{A.5}$$

Let $J_{H_a} = \int (\prod_s F_{H_a,s}) P(\bar{\mathbf{b}} | H_a) d\bar{\mathbf{b}}$; this quantity is also analytically computable by straightforward algebra.

To compute Bayes Factor of H_a versus H_0 under the ES model, we take limits with respect to hyperparameters $(v_1, \dots, v_S, l_1, m_1, \dots, l_S, m_S)$ according to (2.13), that is,

$$\begin{aligned}
\text{BF}^{\text{ES}}(\phi, \omega) &= \lim \frac{\int J_{H_a} \prod_s P(\tau_s) d\tau_1 \cdots d\tau_S}{\int J_{H_0} \prod_s P(\tau_s) d\tau_1 \cdots d\tau_S} \\
&= \frac{\int K_{H_a} d\tau_1 \cdots d\tau_S}{\int K_{H_0} d\tau_1 \cdots d\tau_S}.
\end{aligned} \tag{A.6}$$

Let us denote

$$\text{RSS}_{0,s} = \mathbf{y}_s' \mathbf{y}_s - n_s \bar{y}_s^2, \tag{A.7}$$

$$\text{RSS}_{1,s} = \mathbf{y}_s' \mathbf{y}_s - \mathbf{y}_s' \mathbf{X}_s (\mathbf{X}_s' \mathbf{X}_s)^{-1} \mathbf{X}_s' \mathbf{y}_s, \tag{A.8}$$

$$\delta_s^2 = \frac{1}{\mathbf{g}_s' \mathbf{g}_s - n_s \bar{g}_s^2}, \tag{A.9}$$

$$\hat{\beta}_s = \frac{\mathbf{y}_s' \mathbf{g}_s - n_s \bar{y}_s \bar{g}_s}{\mathbf{g}_s' \mathbf{g}_s - n_s \bar{g}_s^2}, \tag{A.10}$$

$$\zeta^2 = \frac{1}{\sum_s (\delta_s^2 + \phi^2)^{-1}}, \tag{A.11}$$

where \bar{y}_s and \bar{g}_s are the sample means of phenotypes and genotypes in subgroup s . It can be shown that,

$$K_{H_0} = \prod_s \tau_s^{\frac{n_s}{2}-1} \exp\left(-\frac{1}{2} \sum_s \tau_s \cdot \text{RSS}_{0,s}\right), \tag{A.12}$$

and

$$\begin{aligned}
K_{H_a} = & \sqrt{\frac{\zeta^2}{\zeta^2 + \omega^2}} \prod_s \sqrt{\frac{\delta_s^2}{\delta_s^2 + \phi^2}} \\
& \cdot \prod_s \tau_s^{\frac{n_s}{2}-1} \exp \left(-\frac{1}{2} \sum_s \tau_s \left(\frac{\phi^2}{\delta_s^2 + \phi^2} \cdot \text{RSS}_{1,s} + \frac{\delta_s^2}{\delta_s^2 + \phi^2} \cdot \text{RSS}_{0,s} \right) \right) \\
& \cdot \exp \left(\frac{1}{2} \frac{\omega^2 \zeta^2}{\zeta^2 + \omega^2} \left(\sum_s \frac{\hat{\beta}_s \sqrt{\tau_s}}{\delta_s^2 + \phi^2} \right)^2 \right).
\end{aligned} \tag{A.13}$$

The multidimensional integral $\int K_{H_a} d\tau_1 \cdots d\tau_S$ generally does not have a simple analytic form (although it can be represented as finite sums of complicated hypergeometric functions). Next, we show two different approximations, both based on Laplace's method, to evaluate this integral. The first approximation is a direct application of Butler and Wood (2002) and the second one yields a simple analytic expression. Although the integral $\int K_{H_0} d\tau_1 \cdots d\tau_S$ can be analytically computed as a gamma function, for computing the Bayes Factor, we also use Laplace's method to numerically evaluate it (which essentially is applying Sterling's formula) – we find this recipe yields more accurate result for the final Bayes Factor: in particular, when there is only one subgroup ($S = 1$, where the Bayes Factor can be analytically computed as in Servin and Stephens (2007)), we obtain the exact result by applying the first Laplace's approximation.

Laplace's method approximates a multivariate integral in the following way,

$$\int_D h(\boldsymbol{\tau}) e^{g(\boldsymbol{\tau})} d\boldsymbol{\tau} \approx (2\pi)^{S/2} |H_{\hat{\boldsymbol{\tau}}}|^{-1/2} h(\hat{\boldsymbol{\tau}}) e^{g(\hat{\boldsymbol{\tau}})} \tag{A.14}$$

where $\boldsymbol{\tau}$ is an S -vector,

$$\hat{\boldsymbol{\tau}} = \arg \max_{\boldsymbol{\tau}} g(\boldsymbol{\tau}), \tag{A.15}$$

and $|H_{\hat{\boldsymbol{\tau}}}|$ is the absolute value of the determinant of the Hessian matrix of the function g evaluated at $\hat{\boldsymbol{\tau}}$. Note that the factorization of the integrand is rather arbitrary, it only requires that function h is smooth and positively valued and the smooth function g has a unique maximum lying in the interior of

D (for detailed discussion, see Butler (2007)).

Our first approach to apply Laplace's method sets $h(\boldsymbol{\tau}) \equiv 1$. Except for some trivial situations (e.g. $S = 1$), the maximization of $\log K_{H_a}$ with respect to $\boldsymbol{\tau}$ is analytically intractable. In practice, we use the Broyden-Fletcher-Goldfarb-Shanno (BFGS2) algorithm, a gradient-based numerical optimization routine (implemented in the GNU Scientific Library), to perform numerical maximization. This procedure leads to $\widehat{\text{BF}}^{\text{ES}}(\phi, \omega)$.

Alternatively, we apply Laplace's method by factoring the integrand in such a way that g can be analytically maximized. This approach results in a closed-form approximation. More specifically, we factor K_{H_a} into

$$K_{H_a} = h(\tau_1, \dots, \tau_S) e^{g(\tau_1, \dots, \tau_S)}, \quad (\text{A.16})$$

where

$$\begin{aligned} h(\tau_1, \dots, \tau_S) = & \sqrt{\frac{\zeta^2}{\zeta^2 + \omega^2}} \prod_s \sqrt{\frac{\delta_s^2}{\delta_s^2 + \phi^2}} \\ & \cdot \prod_s \exp \left(-\frac{1}{2} \sum_s \frac{\delta_s^2}{\delta_s^2 + \phi^2} \cdot (\text{RSS}_{0,s} - \text{RSS}_{1,s}) \right) \\ & \cdot \exp \left(\frac{1}{2} \frac{\omega^2 \zeta^2}{\zeta^2 + \omega^2} \left(\sum_s \frac{\hat{\beta}_s \sqrt{\tau_s}}{\delta_s^2 + \phi^2} \right)^2 \right) \end{aligned} \quad (\text{A.17})$$

and

$$e^{g(\tau_1, \dots, \tau_S)} = \prod_s \tau_s^{\frac{n_s}{2} - 1} \cdot \exp \left(-\frac{1}{2} \sum_s \tau_s \cdot \text{RSS}_{1,s} \right). \quad (\text{A.18})$$

It is straightforward to show that the unique maximum of $g(\tau_1, \dots, \tau_S)$ is attained at

$$\hat{\tau}_s = \frac{n_s - 2}{\text{RSS}_{1,s}}, \quad s = 1, \dots, S, \quad (\text{A.19})$$

which coincides with the REML estimate of τ_s in subgroup-level regression model (2.1). Following the notations in section 2.3 and noting the relationship between t and F statistics in simple linear regression,

$$T_s^2 = \frac{\text{RSS}_{0,s} - \text{RSS}_{1,s}}{\text{RSS}_{1,s}/(n_s - 2)}. \quad (\text{A.20})$$

Applying (A.14) results in

$$\text{BF}^{\text{ES}}(\phi, \omega) \simeq \sqrt{\frac{\zeta^2}{\zeta^2 + \omega^2}} \exp\left(\frac{\mathcal{T}_{\text{ES}}^2}{2} \frac{\omega^2}{\zeta^2 + \omega^2}\right) \quad (\text{A.21})$$

$$\cdot \prod_s \left(\sqrt{\frac{\delta_s^2}{\delta_s^2 + \phi^2}} \left(\frac{\text{RSS}_{0,s}}{\text{RSS}_{1,s}} \right)^{\frac{n_s}{2}} \exp\left(-\frac{T_s^2}{2} \frac{\delta_s^2}{\delta_s^2 + \phi^2}\right) \right). \quad (\text{A.22})$$

To further simplify the above expression, we use

$$\left(\frac{\text{RSS}_{0,s}}{\text{RSS}_{1,s}} \right)^{n_s/2} = \left(1 + \frac{T_s^2}{n_s - 2} \right)^{\frac{n_s}{2}} = e^{\frac{T_s^2}{2}} \left(1 + O\left(\frac{1}{n_s}\right) \right), \quad (\text{A.23})$$

and (A.21) simplifies to

$$\text{ABF}^{\text{ES}}(\phi, \omega) = \sqrt{\frac{\zeta^2}{\zeta^2 + \omega^2}} \exp\left(\frac{\mathcal{T}_{\text{ES}}^2}{2} \frac{\omega^2}{\zeta^2 + \omega^2}\right) \prod_s \left(\sqrt{\frac{\delta_s^2}{\delta_s^2 + \phi^2}} \exp\left(\frac{T_s^2}{2} \frac{\phi^2}{\delta_s^2 + \phi^2}\right) \right). \quad (\text{A.24})$$

Remarks. Note, in case τ_1, \dots, τ_S are known, we can directly compute the exact Bayes Factor using

$$\text{BF}^{\text{ES}}(\phi, \omega) = \lim \frac{J_{H_a}}{J_{H_0}} \quad (\text{A.25})$$

without evaluating the multi-dimensional integrals in (A.6). In this particular case, it is easy to show that the exact Bayes Factor has the exact functional form as in (2.23) and (A.24), only with all the $\hat{\tau}_s$'s replaced by the corresponding true values of τ_s 's.

Finally, we give the proof for Proposition 1:

Proof of Proposition 1. The derivation above serves as a proof. An alternative proof can be obtained by noting that the REML estimate of $\hat{\tau}$ asymptotically converges to the true value of τ with probability 1. From the remarks above, by applying continuous mapping theorem, we conclude that $\text{ABF}^{\text{ES}}(\phi, \omega)$

converges to $\text{BF}^{\text{ES}}(\phi, \omega)$ with probability 1. □

A.2 Computation in the EE Model

The procedure for computing Bayes Factor assuming an EE model is essentially the same, we omit repeating the details but only show the final results of the Bayes Factor of an EE model H_b , specified by (ψ, w) , versus the null model H_0 ,

$$\text{BF}^{\text{EE}}(\psi, w) = \frac{\int K_{H_b} d\tau_1 \cdots d\tau_S}{\int K_{H_0} d\tau_1 \cdots d\tau_S}. \quad (\text{A.26})$$

The expression of K_{H_0} remains the same as (A.12). We denote

$$\eta^2 = \left(\sum_s \frac{\tau_s}{\delta_s^2 + \tau_s \psi^2} \right)^{-1}. \quad (\text{A.27})$$

It can be shown

$$\begin{aligned} K_{H_b} = & \sqrt{\frac{\eta^2}{\eta^2 + w^2}} \prod_s \sqrt{\frac{\delta_s^2}{\delta_s^2 + \tau_s \psi^2}} \\ & \cdot \prod_s \tau_s^{\frac{n_s}{2} - 1} \exp \left(-\frac{1}{2} \sum_s \tau_s \left(\frac{\tau_s \psi^2}{\delta_s^2 + \tau_s \psi^2} \cdot \text{RSS}_{1,s} + \frac{\delta_s^2}{\delta_s^2 + \tau_s \psi^2} \cdot \text{RSS}_{0,s} \right) \right) \\ & \cdot \exp \left(\frac{1}{2} \frac{w^2}{\eta^2 + w^2} \frac{\left(\sum_s \frac{\tau_s}{\delta_s^2 + \tau_s \psi^2} \cdot \hat{\beta}_s \right)^2}{\eta^2} \right). \end{aligned} \quad (\text{A.28})$$

We use the similar numerical procedure to obtain $\widehat{\text{BF}^{\text{EE}}}(\psi, w)$ as in the ES model.

To derive ABF^{EE} , we factor K_{H_b} into

$$K_{H_b} = h(\tau_1, \dots, \tau_S) e^{g(\tau_1, \dots, \tau_S)}, \quad (\text{A.29})$$

where,

$$\begin{aligned}
h(\tau_1, \dots, \tau_S) = & \sqrt{\frac{\eta^2}{\eta^2 + w^2}} \prod_s \sqrt{\frac{\delta_s^2}{\delta_s^2 + \tau_s \psi^2}} \\
& \cdot \prod_s \exp \left(-\frac{1}{2} \sum_s \frac{\tau_s \delta_s^2}{\delta_s^2 + \tau_s \psi^2} \cdot (\text{RSS}_{0,s} - \text{RSS}_{1,s}) \right) \\
& \cdot \exp \left(\frac{1}{2} \frac{w^2}{\eta^2 + w^2} \frac{\left(\sum_s \frac{\tau_s}{\delta_s^2 + \tau_s \psi^2} \cdot \hat{\beta}_s \right)^2}{\eta^2} \right)
\end{aligned} \tag{A.30}$$

and

$$e^{g(\tau_1, \dots, \tau_S)} = \prod_s \tau_s^{\frac{n_s}{2}-1} \cdot \exp \left(-\frac{1}{2} \sum_s \tau_s \cdot \text{RSS}_{1,s} \right). \tag{A.31}$$

Again, function $g(\tau_1, \dots, \tau_S)$ is maximized at

$$\hat{\tau}_s = \frac{n_s - 2}{\text{RSS}_{1,s}}, \quad s = 1, \dots, S. \tag{A.32}$$

We denote

$$d_s^2 = \frac{1}{\hat{\tau}_s} \delta_s^2 = \frac{\hat{\sigma}_s^2}{\mathbf{g}_s' \mathbf{g}_s - n_s \bar{g}_s^2}, \tag{A.33}$$

$$T_s^2 = \frac{\hat{\beta}_s}{d_s^2}, \tag{A.34}$$

$$\xi^2 = \left(\sum_s \frac{\hat{\tau}_s}{\delta_s^2 + \hat{\tau}_s \psi^2} \right)^{-1} = \frac{1}{\sum_s (d_s^2 + \psi^2)^{-1}}, \tag{A.35}$$

$$\hat{\beta} = \frac{\sum_s (d_s^2 + \psi^2)^{-1} \hat{\beta}_s}{\sum_s (d_s^2 + \psi^2)^{-1}}, \tag{A.36}$$

and

$$\tau_{\text{EE}}^2 = \frac{\hat{\beta}^2}{\xi^2} = \frac{\left(\sum_s \frac{\hat{\beta}}{d_s^2 + \psi^2} \right)^2}{\eta^2}. \tag{A.37}$$

Using the similar procedure as in the ES model, we obtain

$$\text{ABF}^{\text{EE}}(\psi, w) = \sqrt{\frac{\xi^2}{\xi^2 + w^2}} \exp\left(\frac{\tau_{\text{EE}}^2}{2} \frac{w^2}{\xi^2 + w^2}\right) \prod_s \left(\sqrt{\frac{d_s^2}{d_s^2 + \psi^2}} \exp\left(\frac{T_s^2}{2} \frac{\psi^2}{d_s^2 + \psi^2}\right) \right). \quad (\text{A.38})$$

Same as we have discussed in **Remarks** of section A.1, if τ_1, \dots, τ_S are known, the exact Bayes Factor of the EE model has the same function form as in (A.38), only with $\hat{\tau}_s$'s replaced by corresponding τ_s 's.

A.3 Computation using CEFN Priors

Using curved exponential family normal prior, the computation of the Bayes Factors is slightly different than what we show in previous sections. Here, we use the ES model as a demonstration, the procedure for the EE model is very similar.

To compute the Bayes Factor of a CEFN-ES model defined by parameters (k, ω) vs. the null model, we can carry out the same and exact calculation up to (A.5). However, due to the nature of the CEFN prior, we can no longer perform analytic calculation to integrate out \bar{b} . Instead, we exchange the order of integrations by first analytically approximating the multi-dimensional integration with respect to τ_1, \dots, τ_S using the second procedure of Laplace's method described in previous sections. As a result, we obtain the approximate Bayes Factor as a one-dimensional integral

$$\begin{aligned} \text{ABF}_{\text{CEFN}}^{\text{ES}}(k, \omega) &= \frac{1}{\sqrt{2\pi\omega}} \prod_s \left(\frac{\text{RSS}_{0,s}}{\text{RSS}_{1,s}} \right)^{n_s/2} \int_{-\infty}^{\infty} \prod_s \sqrt{\frac{\delta_s^2}{\delta_s^2 + k\bar{b}^2}} \\ &\cdot \exp \left[-\frac{1}{2} \left(\left(\sum_s \frac{1}{\delta_s^2 + k\bar{b}^2} + \frac{1}{\omega^2} \right) \bar{b}^2 - 2 \sum_s \left(\frac{\hat{b}_s}{\delta_s^2 + k\bar{b}^2} \right) \bar{b} + \sum_s \frac{\delta_s^2}{\delta_s^2 + k\bar{b}^2} T_s^2 \right) \right] d\bar{b}. \end{aligned} \quad (\text{A.39})$$

We then apply an adaptive Gaussian quadrature method, QAGI (implemented in GNU scientific library), to numerically evaluate this integral. Essentially, this method first maps the integrand to the semi-open interval $[0, 1)$ using the transformation $y = (1 - \bar{b})/\bar{b}$, then apply the standard adaptive Gaussian quadrature routine for the finite interval integration.

For the EE model with CEFN prior, the final one-dimensional integral can be shown as

$$\begin{aligned} \text{ABF}_{\text{CEFN}}^{\text{EE}}(k, w) &= \frac{1}{\sqrt{2\pi}w} \prod_s \left(\frac{\text{RSS}_{0,s}}{\text{RSS}_{1,s}} \right)^{n_s/2} \int_{-\infty}^{\infty} \prod_s \sqrt{\frac{d_s^2}{d_s^2 + k\bar{\beta}^2}} \\ &\cdot \exp \left[-\frac{1}{2} \left(\left(\sum_s \frac{1}{d_s^2 + k\bar{\beta}^2} + \frac{1}{w^2} \right) \bar{\beta}^2 - 2 \sum_s \left(\frac{\hat{\beta}_s}{d_s^2 + k\bar{\beta}^2} \right) \bar{\beta} + \sum_s \frac{d_s^2}{d_s^2 + k\bar{\beta}^2} T_s^2 \right) \right] d\bar{\beta}. \end{aligned} \quad (\text{A.40})$$

B Bayes Factor for Binary Regression Models

In this section, we show the computation of Bayes Factor for case-control data.

Let us denote $\boldsymbol{\beta}_s = (\mu_s, \beta_s)$. The key component in our computation is to approximate subgroup-level log-likelihood function $l(\boldsymbol{\beta}_s)$ with a quadratic form expanding around its maximum likelihood estimate, i.e.

$$\log P(\mathbf{y}_s | \mathbf{g}_s, \boldsymbol{\beta}_s) = l(\boldsymbol{\beta}_s) \simeq l(\hat{\boldsymbol{\beta}}_s) - \frac{1}{2} (\boldsymbol{\beta}_s - \hat{\boldsymbol{\beta}}_s)' I_s(\hat{\boldsymbol{\beta}}_s) (\boldsymbol{\beta}_s - \hat{\boldsymbol{\beta}}_s), \quad (\text{B.1})$$

where $I_s(\hat{\boldsymbol{\beta}}_s) = \begin{pmatrix} i_{\hat{\mu}_s \hat{\mu}_s} & i_{\hat{\mu}_s \hat{\beta}_s} \\ i_{\hat{\beta}_s \hat{\mu}_s} & i_{\hat{\beta}_s \hat{\beta}_s} \end{pmatrix}$ is the expected Fisher information evaluated at $\hat{\boldsymbol{\beta}}_s$. Although this type of approximation generally requires the observed Fisher information in (B.1), the observed and expected Fisher information indeed coincide as we use the canonical (logistic) link for binary regression model.

Further, we note

$$\gamma_s^2 := \text{Var}(\hat{\beta}_s) = (i_{\hat{\beta}_s \hat{\beta}_s} - i_{\hat{\beta}_s \hat{\mu}_s} i_{\hat{\mu}_s \hat{\mu}_s}^{-1} i_{\hat{\mu}_s \hat{\beta}_s})^{-1} \quad (\text{B.2})$$

is the estimated asymptotic variance of MLE $\hat{\beta}_s$.

Given approximate log-likelihood function (B.1) and a model H_c specified by (ψ, w) , the prior distribution for $\boldsymbol{\beta}_s$ is given by

$$\boldsymbol{\beta}_s | \bar{\boldsymbol{\beta}}, H_c \sim N(\bar{\boldsymbol{\beta}}, \Psi_s), \quad (\text{B.3})$$

where

$$\bar{\beta} = \begin{pmatrix} 0 \\ \bar{\beta} \end{pmatrix} \quad \text{and} \quad \Psi_s = \begin{pmatrix} v_s^2 & 0 \\ 0 & \psi^2 \end{pmatrix}. \quad (\text{B.4})$$

It follows that

$$\begin{aligned} F_{H_c, s} &= \int P(\mathbf{y}_s | \mathbf{g}_s, \beta_s) P(\beta_s | \bar{\beta}, H_c) d\beta_s \\ &= \exp(l(\hat{\beta}_s)) \cdot |\Psi_s|^{-\frac{1}{2}} \cdot |I_s + \Psi_s^{-1}|^{-\frac{1}{2}} \cdot \exp\left(-\frac{1}{2} \hat{\beta}_s' (I_s - I_s(I_s + \Psi_s^{-1})^{-1} I_s) \hat{\beta}_s\right) \\ &\quad \cdot \exp\left(-\frac{1}{2} \left(\bar{\beta}' (\Psi_s^{-1} - \Psi_s^{-1} (I_s + \Psi_s^{-1})^{-1} \Psi_s^{-1}) \bar{\beta} - \bar{\beta}' \eta_s - \eta_s' \bar{\beta} \right)\right), \end{aligned} \quad (\text{B.5})$$

with $\eta_s = \Psi^{-1} (I_s + \Psi^{-1})^{-1} I_s \hat{\beta}_s$.

Under contrasting null model H_0 , the parameter space is restricted to $\beta_s = 0$, for β_s satisfies this restriction

$$(\beta_s - \hat{\beta}_s)' I_s (\hat{\beta}_s) (\beta_s - \hat{\beta}_s) = i_{\hat{\mu}_s \hat{\mu}_s} \cdot (\mu_s - \hat{m}_s)^2 + \frac{\hat{\beta}_s^2}{\gamma_s^2}, \quad (\text{B.6})$$

where $\hat{m}_s = \hat{\mu}_s + \frac{i_{\hat{\mu}_s \hat{\beta}_s}}{i_{\hat{\mu}_s \hat{\mu}_s}} \hat{\beta}_s$. It can be shown that

$$\begin{aligned} F_{H_0, s} &= \int P(\mathbf{y}_s | \mathbf{g}_s, \beta_s) P(\beta_s | \bar{\beta}, H_0) d\beta_s \\ &= \exp(l(\hat{\beta}_s)) \cdot v_s^{-1} (i_{\hat{\mu}\hat{\mu}} + v_s^{-2})^{-\frac{1}{2}} \cdot \exp\left(-\frac{\hat{\beta}_s^2}{2\gamma_s^2}\right) \\ &\quad \cdot \exp\left(-\frac{1}{2} \left(\hat{m}_s' i_{\hat{\mu}\hat{\mu}} \hat{m}_s - (i_{\hat{\mu}\hat{\mu}} \hat{m}_s)' (i_{\hat{\mu}\hat{\mu}} + v_s^{-2})^{-1} (i_{\hat{\mu}\hat{\mu}} \hat{m}_s) \right)\right). \end{aligned} \quad (\text{B.7})$$

Finally, we compute

$$\text{ABF}^{\text{CC}}(\psi, w) = \lim \frac{\int (\prod_s F_{H_c, s}) P(\bar{\beta} | H_c) d\bar{\beta}}{\prod_s F_{H_0, s}}, \quad (\text{B.8})$$

where the limit is taken as $v_s \rightarrow \infty, \forall s$. By straightforward algebra, we obtain the following final result

$$\text{BF}^{\text{CC}}(\psi, w) \approx \text{ABF}^{\text{CC}}(\psi, w) := \text{ABF}_{\text{single}}^{\text{CC}}(\mathcal{Z}_{\text{cc}}^2, \xi; w) \cdot \prod_s \text{ABF}_{\text{single}}^{\text{cc}}(Z_s^2, \gamma_s; \psi), \quad (\text{B.9})$$

where

$$\text{ABF}_{\text{single}}^{\text{CC}}(Z_s^2, \gamma_s; \psi) = \sqrt{\frac{\gamma_s^2}{\gamma_s^2 + \psi^2}} \exp\left(\frac{Z_s^2}{2} \frac{\psi^2}{\gamma_s^2 + \psi^2}\right), \quad (\text{B.10})$$

$$\text{ABF}_{\text{single}}^{\text{CC}}(\mathcal{Z}_{\text{cc}}^2, \xi; w) = \sqrt{\frac{\xi^2}{\xi^2 + w^2}} \exp\left(\frac{\mathcal{Z}_{\text{cc}}^2}{2} \frac{w^2}{\xi^2 + w^2}\right), \quad (\text{B.11})$$

$$(\text{B.12})$$

and

$$\gamma_s^2 := \text{se}(\hat{\beta}_s)^2, \quad (\text{B.13})$$

$$Z_s^2 = \frac{\hat{\beta}_s^2}{\gamma_s^2}, \quad (\text{B.14})$$

$$\hat{\beta} = \frac{\sum_s (\gamma_s^2 + \psi^2)^{-1} \hat{\beta}_s}{\sum_s (\gamma_s^2 + \psi^2)^{-1}}, \quad (\text{B.15})$$

$$\xi^2 := \text{se}(\hat{\beta})^2 = \frac{1}{\sum_s (\gamma_s^2 + \psi^2)^{-1}}, \quad (\text{B.16})$$

$$\mathcal{Z}_{\text{cc}}^2 = \frac{\hat{\beta}^2}{\xi^2}. \quad (\text{B.17})$$

C Small Sample Size Correction for Approximate Bayes Factors

The accuracy of ABF^{ES} relies on the sample sizes in subgroups: when sample sizes are small in some subgroups, the approximation may become inaccurate. In particular, we consider the behavior of the approximate Bayes Factor when the null hypothesis is true. A valid Bayes Factor has the property that

$$\text{E}(\text{BF}|H_0) = 1, \quad (\text{C.1})$$

where the expectation is taken with respect to the data distribution (\mathbf{Y} in our settings) under the null model. This is because,

$$\mathbb{E}(\text{BF}|H_0) = \int \frac{P(\mathbf{Y}|H_1)}{P(\mathbf{Y}|H_0)} \cdot P(\mathbf{Y}|H_0) d\mathbf{Y} = 1. \quad (\text{C.2})$$

Unfortunately, when sample sizes are small, (C.1) can be violated (as the expected value is strictly greater than 1) and this essentially indicates that the approximation becomes inaccurate.

To demonstrate a violation of (C.1), we consider the special case of one single subgroup. The approximate Bayes Factor assuming the ES model with parameters (ϕ, ω) is given by

$$\text{ABF}_{\text{single}}^{\text{ES}}(\phi, \omega) = \sqrt{1 - \lambda} \exp\left(\frac{\lambda}{2} T_s^2\right), \quad (\text{C.3})$$

and,

$$\log(\text{ABF}_{\text{single}}^{\text{ES}}(\phi, \omega)) = \frac{1}{2} \log(1 - \lambda) + \frac{\lambda}{2} T_s^2, \quad (\text{C.4})$$

where $\lambda = \frac{\phi^2 + \omega^2}{\phi^2 + \omega^2 + \delta_s^2}$ and takes values from $[0, 1]$. Under H_0 , T_s follows t-distribution with $n_s - 2$ degree of freedom and

$$\mathbb{E}(T_s^2|H_0) = \frac{n_s - 2}{n_s - 4} > 1. \quad (\text{C.5})$$

Now consider the continuous function

$$f(\lambda) = \frac{1}{\lambda} \log\left(\frac{1}{1 - \lambda}\right) \quad (\text{C.6})$$

for $\lambda \in [0, 1]$, it can be shown that

$$\lim_{\lambda \rightarrow 0} f(\lambda) = 1 \quad (\text{C.7})$$

$$\lim_{\lambda \rightarrow 1} f(\lambda) = \infty. \quad (\text{C.8})$$

Hence, there must exist values of $\lambda \in (0, 1)$, such that

$$1 < f(\lambda) < E(T_s^2 | H_0). \quad (\text{C.9})$$

Consequently, by Jensen's inequality, for those λ values

$$\log(E(\text{ABF}_{\text{single}}^{\text{ES}} | H_0)) \geq E(\log(\text{ABF}_{\text{single}}^{\text{ES}}) | H_0) > 0. \quad (\text{C.10})$$

This shows property (C.1) does not generally hold for approximate Bayes Factors, and when sample size n_s is small, the inaccuracy may become severe.

We now propose a simple correction procedure for small sample sizes, which ensures the resulting approximation satisfies property (C.1). Specifically, we modify (2.23) into the following form

$$\text{A}^* \text{BF}^{\text{ES}}(\phi, \omega) = \text{ABF}_{\text{single}}^{\text{ES}}(q(\mathcal{T}_{\text{ES}}^2), \zeta; \omega) \cdot \prod_s \text{ABF}_{\text{single}}^{\text{ES}}(q_s(T_s^2), \delta_s; \phi), \quad (\text{C.11})$$

where the function q_s denotes a one-to-one quantile transformation from a t-distribution with $n_s - 2$ degree of freedom to a standard normal distribution, and the function q is defined as

$$q(\mathcal{T}_{\text{ES}})^2 = \frac{\hat{b}_{\text{cor}}^2}{\zeta^2}, \quad (\text{C.12})$$

where

$$\hat{b}_{\text{cor}} = \frac{\sum_s (\delta_s^2 + \phi^2)^{-1} \delta_s q_s(T_s)}{\sum_s (\delta_s^2 + \phi^2)^{-1}}. \quad (\text{C.13})$$

Note, the quantile transformation functions q_s and q converge to the identity mappings as $n_s \rightarrow \infty$ and the asymptotic property of (2.23) is preserved. The numerical performance of this correction is demonstrated in appendix D.

To show the corrected version of approximate Bayes Factor satisfying (C.1), we note that ABF^{ES} depends on data \mathbf{Y} only through T_s (δ_s depends on genotype data but not \mathbf{Y}). Further, from **Remark** in appendix A.1, we also notice the approximation becomes an exact Bayes Factor (for which property

(C.1) is guaranteed) if estimated error variance terms $\hat{\sigma}_s^2$'s are replaced by their corresponding true values. When the true error variances are plugged in, under the H_0 , T_s 's instead follow standard normal distributions. It is therefore sufficient to satisfy property (C.1) by quantile transforming each individual T_s in (2.23) from the t-distribution to standard normal distribution. In essence, the correction can be viewed as a general strategy of providing a better point estimate of residual errors, therefore the similar strategy also likely improves the accuracy of approximate Bayes Factor when EE model or CEFN model is used.

D Numerical Accuracy of Bayes Factor Evaluations

In this section, we evaluate the numerical accuracy of various approximation methods for computing the Bayes Factors.

We use the dataset from population eQTL study (Stranger et al. (2007)) discussed in section 3.3 for this purpose. For each of the 8,427 genes examined, we select the top associated *cis*-SNP based on the values of $\widehat{\text{BF}}_{\text{av}}^{\text{ES}}$ and re-calculate the Bayes Factor directly based on (A.6) using a general adaptive Gaussian quadrature procedure (Note, because of its high computational cost in numerically evaluating multi-dimensional integrals, this numerical recipe does not apply in general practice). We treat these results as the “truth” and make comparison with $\widehat{\text{BF}}_{\text{av}}^{\text{ES}}$ and $\text{ABF}_{\text{av}}^{\text{ES}}$ (with and without small sample corrections). Moreover, we convert various numerical results of Bayes Factors to log 10 scale and compute Root Mean Squared Errors (RMSE) for each approximation.

The results of the numerical evaluation for the ES model shown in Table 5 and Figure 6. Although the sample sizes in each subgroup are quite small in this dataset (41 Europeans, 59 Asians and 41 Africans), the numerical results of $\widehat{\text{BF}}_{\text{av}}^{\text{ES}}$ are almost identical to the results obtained from the adaptive Gaussian quadrature procedure (RMSE = 1.2×10^{-4} in log 10 scale). As expected, the approximate Bayes Factor, $\text{ABF}_{\text{av}}^{\text{ES}}$, has the worst numerical performance, mainly due to the small sample sizes in this dataset. Nevertheless, the ranking of the SNPs by $\text{ABF}_{\text{av}}^{\text{ES}}$ is quite consistent with true values (rank correlation = 0.99). Figure 6 suggests that under small sample situations, $\text{ABF}_{\text{av}}^{\text{ES}}$ tends to over-evaluate the true value and this over-evaluation can become quite severe when the true values are extremely large.

On the other hand, the proposed small sample size correction method seems very effective: with this simple correction, the resulting $A^*BF_{av}^{ES}$ are quite accurate comparing with the true values.

We also perform a similar experiment for the EE model using the same dataset with five levels of $\sqrt{\psi^2 + w^2}$ values: 0.1, 0.2, 0.4, 0.8, 1.6, and seven degrees of heterogeneities characterized by ψ^2/w^2 values: 0, 1/4, 1/2, 1, 2, 4, ∞ , and we assign these 35 grid values equal prior weight. The results are similar with the case in the EE model and shown in Table 6.

	$\log_{10}(\widehat{BF}_{av}^{ES})$	$\log_{10}(ABF_{av}^{ES})$	$\log_{10}(A^*BF_{av}^{ES})$
RMSE	1.2×10^{-4}	4.95	0.14

Table 5: Numerical accuracy of three approximations for evaluating Bayes Factors under the ES model. \widehat{BF}_{av}^{ES} is based on the first approximation of Laplace's method discussed in appendix A, ABF_{av}^{ES} is computed using (2.23) and $A^*BF_{av}^{ES}$ is based on (C.11) which is corrected for small sample sizes.

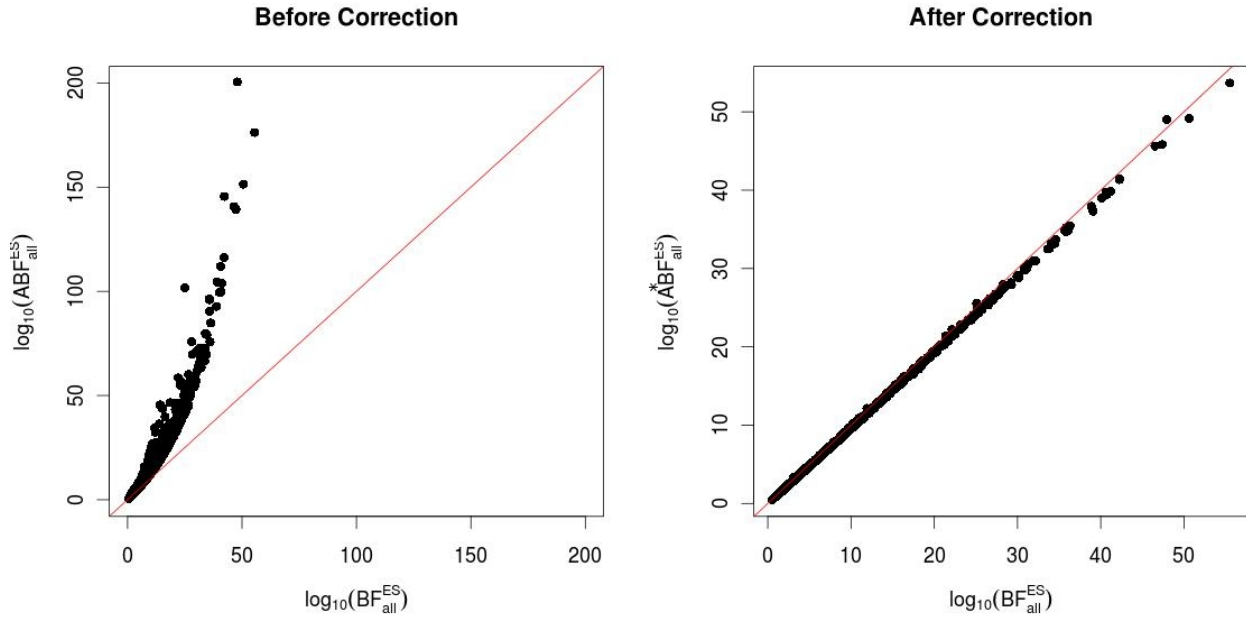


Figure 6: Comparison of approximate Bayes Factors before and after applying small sample size corrections.

	$\log_{10}(\widehat{\text{BF}}_{\text{av}}^{\text{EE}})$	$\log_{10}(\text{ABF}_{\text{av}}^{\text{EE}})$	$\log_{10}(\text{A}^*\text{BF}_{\text{av}}^{\text{EE}})$
RMSE	4.1×10^{-4}	5.03	0.09

Table 6: Numerical accuracy of three approximations for evaluating Bayes Factors under the EE model.

References

- D. M. Bravata and I. Olkin. Simple Pooling versus Combining in Meta-Analysis. *Evaluation & the Health Professions*, 24(2):218–230, June 2001.
- Stephen Burgess, Simon G Thompson, S Burgess, S G Thompson, G Andrews, et al. Bayesian methods for meta-analysis of causal relationships estimated using genetic instrumental variables. *Statistics in medicine*, 29(12):1298–311, May 2010.
- R. Butler. *Saddlepoint Approximations with Applications*. Cambridge University Press, 1st edition, 2007.
- Roland W. Butler and Andrew T. A. Wood. Laplace approximations for hypergeometric functions with matrix argument. *The Annals of Statistics*, 30(4):1155–1177, August 2002.
- Maria De Iorio, Paul J Newcombe, Ioanna Tachmazidou, Claudio J Verzilli, and John C Whittaker. Bayesian semiparametric meta-analysis for genetic association studies. *Genetic Epidemiology*, 35(5):333–340, 2011.
- Antigone S Dimas, Samuel Deutsch, Barbara E Stranger, Stephen B Montgomery, Christelle Borel, et al. Common regulatory variation impacts gene expression in a cell type-dependent manner. *Science (New York, N.Y.)*, 325(5945):1246–50, September 2009.
- WH DuMouchel and GE Harris. Bayes methods for combining the results of cancer studies in humans and other species with discussion. *Journal of the American Statistical Association*, 78:293–315, 1983.
- Richard M. Durbin, David L. Altshuler, Gonçalo R. Abecasis, David R. Bentley, Aravinda Chakravarti, et al. A map of human genome variation from population-scale sequencing. *Nature*, 467(7319):1061–1073, October 2010.

- DM Eddy, V Hasselblad, and R Schachter. A bayesian method for synthesizing evidence. *International Journal of Technical Assistance in Health Care*, 6:31–55, 1990.
- Adi Fledel-Alon, Ellen Miranda Leffler, Yongtao Guan, Matthew Stephens, Graham Coop, et al. Variation in human recombination rates and its genetic determinants. *PloS one*, 6(6):e20321, January 2011.
- GH Givens, DD Smith, and RL Tweedie. Bayesian data-augmented meta-analysis that account for publication bias issues exemplified in the passive smoking debate. *Statistical Science*, 12:221–250, 1997.
- Buhm Han and Eleazar Eskin. Random-Effects Model Aimed at Discovering Associations in Meta-Analysis of Genome-wide Association Studies. *The American Journal of Human Genetics*, 88(5):586–598, 2011.
- Valen E Johnson. Bayes factors based on test statistics. *Journal of the Royal Statistical Society - Series B: Statistical Methodology*, 67(5):689–701, 2005.
- Valen E Johnson. Properties of Bayes Factors Based on Test Statistics. *Scandinavian Journal of Statistics*, 35(2):354–368, 2008.
- Augustine Kong, Gudmar Thorleifsson, Hreinn Stefansson, Gisli Masson, Agnar Helgason, et al. Sequence variants in the RNF212 gene associate with genome-wide recombination rate. *Science (New York, N.Y.)*, 319(5868):1398–401, March 2008.
- Sridhar Kudaravalli, Jean-Baptiste Veyrieras, Barbara E Stranger, Emmanouil T Dermitzakis, and Jonathan K Pritchard. Gene expression levels are a target of recent natural selection in the human genome. *Molecular biology and evolution*, 26(3):649–58, March 2009.
- Jeremie J Lebrech, Theo Stijnen, and Hans C Van Houwelingen. Statistical Applications in Genetics and Molecular Biology Dealing with Heterogeneity between Cohorts in Genomewide SNP Association Studies Dealing with Heterogeneity between Cohorts in Genomewide SNP Association Studies. *Statistical Applications in Genetics and Molecular Biology*, 9(1), 2010.

- Z Li and CB Begg. Random effects models for combining results from controlled and uncontrolled studies in a meta-analysis. *Journal of the American Statistical Association*, 89:1523–1527, 1994.
- A L Mila and H K Ngugi. A Bayesian approach to meta-analysis of plant pathology studies. *Phytopathology*, 101(1):42–51, January 2011.
- Carole Ober, Dagan A Loisel, and Yoav Gilad. Sex-specific genetic architecture of human disease. *Nature reviews. Genetics*, 9(12):911–22, December 2008.
- Art B. Owen. Karl Pearsons meta-analysis revisited. *The Annals of Statistics*, 37(6B):3867–3892, December 2009.
- Joseph K Pickrell, John C Marioni, Athma A Pai, Jacob F Degner, Barbara E Engelhardt, et al. Understanding mechanisms underlying human gene expression variation with RNA sequencing. *Nature*, 464(7289):768–72, April 2010.
- Bertrand Servin and Matthew Stephens. Imputation-based analysis of association studies: candidate regions and quantitative traits. *PLoS genetics*, 3(7):e114, July 2007.
- Dalene K Stangl and Donald A Berry. *Meta-Analysis in Medicine and Health Policy*. Marcel Dekker, Inc, 2000.
- Barbara E Stranger, Alexandra C Nica, Matthew S Forrest, Antigone Dimas, Christine P Bird, et al. Population genomics of human gene expression. *Nature genetics*, 39(10):1217–24, October 2007.
- A J Sutton and K R Abrams. Bayesian methods in meta-analysis and evidence synthesis. *Statistical methods in medical research*, 10(4):277–303, August 2001.
- Tanya M. Teslovich, Kiran Musunuru, Albert V. Smith, Andrew C. Edmondson, Ioannis M. Stylianou, et al. Biological, clinical and population relevance of 95 loci for blood lipids. *Nature*, 466(7307):707–713, August 2010.
- The International HapMap Consortium. A haplotype map of the human genome. *Nature*, 437:1299–1320, 2005.

- Claudio Verzilli, Tina Shah, Juan P Casas, Juliet Chapman, Manjinder Sandhu, et al. Bayesian Meta-Analysis of Genetic Association Studies with Different Sets of Markers. *The American Journal of Human Genetics*, 82(April):859–872, 2008.
- Jean-Baptiste Veyrieras, Sridhar Kudaravalli, Su Yeon Kim, Emmanouil T Dermizakis, Yoav Gilad, et al. High-resolution mapping of expression-QTLs yields insight into human gene regulation. *PLoS genetics*, 4(10):e1000214, October 2008.
- Jon Wakefield. Bayes factors for genome-wide association studies: comparison with P-values. *Genetic epidemiology*, 33(1):79–86, January 2009.
- A Whitehead and J Whitehead. A general parametric approach to the meta-analysis of randomized clinical trials. *Statistics in Medicine*, 10:1665–1677, 1991.
- Cristen J Willer, Yun Li, and Gonçalo R Abecasis. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics (Oxford, England)*, 26(17):2190–1, September 2010.
- Eleftheria Zeggini, Laura J Scott, Richa Saxena, Benjamin F Voight, Jonathan L Marchini, et al. Meta-analysis of genome-wide association data and large-scale replication identifies additional susceptibility loci for type 2 diabetes. *Nature Genetics*, 40(5):638–645, March 2008.